

HIGH-PERFORMANCE PARALLEL INTERFACE - Scheduled Transfer (HIPPI-ST)

March 31, 1997

Secretariat:

Information Technology Industry Council (ITI)

ABSTRACT: This standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

NOTE:

This is an internal working document of X3T11, a Technical Committee of Accredited Standards Committee X3. As such, this is not a completed standard. The contents are actively being modified by X3T11. This document is made available for review and comment only. For current information on the status of this document contact the individuals shown below:

POINTS OF CONTACT:

Roger Cummings (X3T11 Chairman)
Distributed Processing Technology
140 Candace Drive
Maitland, FL 32751
(407) 830-5522 x348, Fax: (407) 260-5366
E-mail: cummings_roger@dpt.com

Ed Grivna (X3T11 Vice-Chairman)
Cypress Semiconductor
2401 East 86th Street
Bloomington, MN 55425
(612) 851-5200, Fax: (612) 851-5087
E-mail: elg@cypress.com

Don Tolmie (HIPPI-ST Technical Editor)
Los Alamos National Laboratory
CIC-5, MS-B255
Los Alamos, NM 87545
(505) 667-5502, Fax: (505) 665-7793
E-mail: det@lanl.gov

Comments on Rev 0.5

This is a preliminary document undergoing lots of changes. Many of the additions are just place holders, or are put there to stimulate discussion. Hence, do not assume that the items herein are correct, or final – everything is subject to change. This page tries to outline where we are; what has been discussed and semi-approved, and what has been added or changed recently and deserves your special attention. This summary relates to changes since the previous revision. Also, previous open issues are outlined with a single box, new open issues ones are marked with a double bar on the left edge of the box.

Changes are marked with margin bars so that changed paragraphs are easily found, and then highlights mark the specific changes. The list below just describes the major changes, for detail changes please compare this revision to the previous revision.

Please help us in this development process by sending comments, corrections, and suggestions to the Technical Editor, Don Tolmie, of the Los Alamos National Laboratory, at det@lanl.gov. If you would like to address the whole group working on this document, send the comment(s) to hippie@network.com.

1. Did a global change of "setup" to "set up".
2. In 1, Scope, changed the 5th bullet from "...maximum size that will not..." to "...maximum size transmission unit that will not...".
3. In 2, added a normative reference to HIPPI-SC; now needed for annex B (mapping to HIPPI-FP).
4. In 3.1, added definitions for Data Operation and log.
5. In 3.2, changed the example from "DATA" to "END" since we are going away from Data being all capitals.
6. In 4.2, changed "...of an STU, together totaling up to no more than 4 gigabytes (2^{32} bytes)" to "...of an STU of up to 2 gigabytes (2^{31} bytes)".
7. In figure 3, changed the size of the data payload from "Total $\leq (2^{32} - 64)$ bytes (about 4 gigabytes)" to " $\leq 2^{31}$ bytes".
8. In 4.2, next to last bullet, changed "...Retry..." to "...Retry and Timeout...".
9. In figure 4, changed the presentation of the three items in the upper left corner to be more of a tuple style. Changed "STU size" to "Max-STU Size". Added "T_len" to the Transfer Descriptor. Changed the arrows between the Block Descriptor and Buffer Descriptor Table to angled lines.
10. In 4.2, next to last paragraph, changed "...data size (in bytes)..." to "...Transfer length (T_len, in bytes)...", and changed "...identifies..." to "...identify...".
11. In 4.3, changed "...and allocate..." to "...by allocating...".
12. In 4.3.2, added at the end: "For example, if the HIPPI-ST is used to encapsulate TCP/IP, then the EtherType would be x'0800'. If HIPPI-ST is being used to encapsulate legacy HIPPI-FP or user data, then the EtherTypes would be x'8180' and x'8181' respectively. If the incoming Port number is invalid, then the Operations shall not be executed (see 9.4.1).
13. In 4.3.3, added "(see 9.4.1)" at the very end.
14. In 4.3.4, changed "...an integral power of two" to "...an integral power of two, i.e., 2^x where $8 \leq x \leq 32$ ".
15. In 4.3.5, changed "...STU" to "...STU (see 4.4.9)"; changed "...declares its..." to "...declares the..."; and changed "...which must be..." to "The Max-STU size must be...". In the 2nd paragraph, changed "...an integral power of two" to "...an integral power of two, i.e., 2^x where $8 \leq x \leq 32$ ".
16. In 4.3.6, changed "Slots" to "Slot" in several places. Changed "...Notify or Interrupt flag is set" to "...Silent = 0 or Interrupt = 1". Did a global change of "DATA Operation" to "Data Operation". In the next to last paragraph, changed "...thus adjusting..." to "...which adjusts...".
17. In 4.4.3, changed "...shall mean..." to "...shall indicate...".
18. In 4.4.2, changed the nomenclature so that rather than setting up a Transfer between G and H, it is now S and R for consistency within table 3.
19. In 4.4.5, changed "Block size" to "Blocksize", i.e., all one word. In the first sentence,

- changed "...Block size for..." to "...Blocksize (in bytes) for...".
20. In 4.4.6, a major change is that the STU counter now starts at zero and increments; before it started at the max value and decremented to zero. Added that the last STU of a Block is marked with Last = 1. Changed the last sentence from "...buffer region boundary, Block size..." to "...buffer boundary, Blocksize...".
 21. In 4.4.7, changed what had been a note into text describing the initial offset, and rewrote the text (3rd paragraph). Added the 4th paragraph stating that best performance is achieved when using Offset = 0.
 22. In 4.4.9, changed "The last example..." to "Example (c)..." . Added the summary of size relationships.
 23. In 4.5.1, changed "This mechanism..." to "Flow control...".
 24. In 4.5.2, added "i.e., to match State_Request and State_Response Operations" at the end of the paragraph.
 25. In 4.5.3, changed "...response (with Reject = 1) when..." to "...response when...".
 26. In 4.5.4, changed "...occurring..." to "...which occur...", and changed "...media..." to "...medium...".
 27. In 6, changed from "The Schedule Header shall be used with all Scheduled Transfer Operations" to "The Schedule Header fields are named for the most common parameter for which the field is used. Many of the fields have different uses depending on the Operation type, and some Operations do not use one or more of the fields at all." . Changed "...parameter..." to "...field...".
 28. In 6.1, changed the title from "...parameters" to "...fields". Checked the usage of these words throughout the document and fixed the words several places, e.g., three places in the first paragraph of 6.1.
 29. In 6.1, under the S_count item, changed the 3rd bullet from "DATA" to "Data", and changed "...the STU counter, counting down to zero..." to "...the STU number...".
 30. In 6.1, under the Offset item, 3rd bullet, changed "...highest number..." to "...highest numbered...".
 31. In 6.2, under the Notify flag, changed the whole paragraph to "Silent". Changed figure 8 from "Notify" and "N" to "Silent" and "T". Changed the text in Interrupt, and in the note following to reflect the change from Notify to Silent.
 32. In 6.2, under the Send_State flag, changed "...Interrupt or Notify flag must also = 1" to "...Interrupt = 1 or Silent = 0 must be true".
 33. In figure 8, changed the "First" (F) flag to "Last" (L).
 34. In 6.2, under the Interrupt flag, changed "Requests exactly the same actions as the Notify flag (see above) on both DATA and non-DATA Operations with the additional requirement that a signal be delivered to the upper-layer entity (see 4.5.5)." to "Requests that a signal, and the Schedule Header with Interrupt = 1, be delivered to the appropriate upper-layer entity after delivery of any data payload to the same entity. Interrupt signals shall be delivered for Control Operations only when Interrupt = 1 (i.e., unlike Notify, Interrupts are not implied for Control Operations). (See 4.5.5.)".
 35. In 6.2, changed the "First" flag to "Last", and changed "The first STU..." to "The last STU...".
 36. In 6.2 under Data Channel assignment, changed "...and is copied in each DATA Operation" to "...and is the Data Channel to be used for Data Operations associated with this Transfer." . Changed "...(2^{17} - 64) bytes, i.e., approximately..." to "... 2^{17} bytes, i.e., 128..." , and changed "...(2^{32} - 64) bytes, i.e., approximately 4 gigabytes" to "... 2^{31} bytes (i.e., 2 gigabytes)". Note that the max size changed. Changed the note from "The maximums were sized to accommodate the Schedule Header and 24 bytes of lower-layer protocol header" to "Data Channel assignment value b'00' is reserved".
 37. In 7.1, under EtherType, changed "...characterizes the data payloads..." to "...characterizes the ULP data payloads...".
 38. In 7.2, under B-Max-STU, changed "...B-Max-Size value..." to "...B-Max-STU value...".
 39. In 8.1, changed "...shall be allocated for a Transfer" to "...shall be allocated, and authorization be given, for a Transfer".

- Under T_len, added "or that the size of the Transfer is unlimited".
40. In 8.2, capitalized "Transfer" in the first sentence.
 41. In 8.3, added "issued by the Final Destination" in the first sentence. Under Offset, changed "...wants to..." to "...must...". Added the R_id parameter to the semantics, and to the text.
 42. In 8.4, under Offset, changed "...wants to..." to "...must...".
 43. In 8.5, changed "DATA" to "Data" globally. Added "Opaque [OS_Bufx], Opaque [OS_Offset]" to the semantics list. Under Flags, changed "Notify" to "Silent", and deleted "...and either Interrupt or Notify = 1". Under S_count, deleted the sentence reading "S_count shall be decremented before each STU is sent". Added the paragraph describing the opaque data.
 44. In 8.6, under Issued, changed "...a State_Responses that it expected..." to "...the State_Response that it expected from a Data Operation...".
 45. In table 3, changed the title from "...G and H" to "...S and R". Changed the G_id and H_id values to S_id and R_id. Added the "T_len" parameter for CTS. Added the Interrupt flag to CTS and SR. Changed the "F" bit to "L", "N" to "T", and the "D" bit to non-bold, non-italic. In note 2, changed "F = First..." to "L = Last...", and "N = Notify" to "T = Silent".
 46. In 9, added the last two sentences in the paragraph.
 47. In 9.1, changed "...media discards or damages..." to "media discard or damage...". Did a global change of "Op-timeout" to "Op_timeout", i.e., underscore instead of hyphen. Added the sentence "Op_timeout_Occurances shall be logged.". In the last paragraph, changed "...system dependent..." to "...system and/or Port dependent...". Added the sentence reading "Max_Retry_Occurances shall be logged..".
 48. In 9.2, third paragraph, changed "...these data movement Operations" to "...these Operations". Added the sentence reading "The ULP may or may not use Op_timeout to indicate failure..".
 49. In table 4, deleted "State_Response, or Request_To_Receive when Persistent = 1" as responses to Request_To_Send.
 50. In 9.3.2, changed "...Request_To_Send_Response)..." to "...Request_Port)...".
 51. In 9.4.1, deleted the sentence reading "All Operations, excluding Request_Port, should have a Source Port value that matches the state for this Virtual Connection.". Changed "...invalid Port..." to "...invalid Destination Port value...".
 52. In 9.4.4, changed "...less than 256 bytes, or is not a power of 2..." to "...less than 256 bytes..." in two places in the paragraph.
 53. In 9.5.3, changed "...shall be discarded..." to "...shall be rejected (see 4.5.3)...".
 54. In 9.5.6, last paragraph, changed "...Operation without the "First STU of Block" flag bit contains an S_count that is not one less than..." to "Operation contains and S_count that is not one greater than...", and changed "... (for this Scheduled Transfer)..." to "... (for this Block)...".
 55. In table 5, added entries for "Max_Retry_Occurance" and "Op_timeout_Occurance". Added abbreviations for PT, PTA, and RTS.
 56. In annex A, changed the title to "Using lower layer protocols" with subsections for HIPPI-6400-PH and HIPPI-FP. If we add more lower layer protocols, then it won't require changing the annex numbers.
 57. In annex A.1, changed "...shall use Data Channel 1, 2, or 3..." to "...shall use Virtual Channel 1, 2, or 3...". Deleted the paragraph about HIPPI-ST providing the pad if the payload size is not a multiple of 32 bytes. Added a paragraph, and bullets, detailing the values of M_len for Control and Data Operations.
 58. Annex A.2 is extracted from Craig Davidson's earlier proposal, and needs to be checked.
 59. Figure A.1 was changed from "(in all except DATA Operations)" to "(in Control Operations)". The maximum size was also changed from " 2^{32} - 64) bytes (approx 4 gigabytes)" to " 2^{31} bytes (2 gigabytes) of".
 60. Added all of annex B about striping.

Contents

	Page
Foreword	ix
Introduction	x
1 Scope	1
2 Normative references.....	1
3 Definitions and conventions	2
3.1 Definitions	2
3.2 Editorial conventions	3
3.2.1 Binary notation	3
3.2.2 Hexadecimal notation	3
3.3 Acronyms and other abbreviations	3
4 System overview	3
4.1 Control Channels and Data Channels	3
4.2 System model	4
4.3 Virtual Connections	6
4.3.1 Sequences and Operations.....	6
4.3.2 Ports.....	6
4.3.3 Keys	7
4.3.4 Buffer size (Bufsize)	7
4.3.5 Max-STU size	7
4.3.6 Slots and Sync parameter.....	7
4.3.7 Concatenate	8
4.3.8 Source_Concatenate	8
4.3.9 Persistent	8
4.4 Data movement.....	9
4.4.1 Sequences and Operations.....	9
4.4.2 Transfer identifiers (R_id and S_id)	9
4.4.3 Transfer length (T_len)	9
4.4.4 Blocks.....	9
4.4.5 Blocksize	10
4.4.6 STUs	10
4.4.7 Bufx and Offset	10
4.4.8 OS_Bufx and OS_Offset	10
4.4.9 Packing examples	10
4.5 Operations management.....	11
4.5.1 Flow control	11
4.5.2 Status Operations	12
4.5.3 Rejected Operations	12
4.5.4 Lost Operations	12
4.5.5 Interrupts	12
5 Service interface.....	13
5.1 Service primitives.....	13
5.2 Sequences of primitives	13
6 Schedule Header	14
6.1 Schedule Header fields	14
6.2 Scheduled Transfer flags.....	15

7	Virtual Connection management	16
7.1	Request_Port	16
7.2	Request_Port_Response	17
7.3	Port_Teardown	17
7.4	Port_Teardown_ACK	18
7.5	Port_Teardown_Complete	18
8	Data movement	19
8.1	Request_To_Send	19
8.2	Request_To_Send_Response	19
8.3	Request_To_Receive	20
8.4	Clear_To_Send	21
8.5	Data	21
8.6	Request_State	22
8.7	State_Response	23
8.8	END	24
8.9	END_ACK	25
9	Error processing	26
9.1	Operation timeout	26
9.2	Operation Pairs	26
9.3	Syntax errors	26
9.3.1	Undefined Opcode	26
9.3.2	Unexpected Opcode	26
9.4	Virtual Connection errors	26
9.4.1	Invalid Key or Port	26
9.4.2	Slots exceeded	27
9.4.3	Unknown EtherType	27
9.4.4	Illegal Bufsize	27
9.4.5	Illegal STU size	27
9.5	Scheduled Transfer errors	27
9.5.1	Invalid S_id	27
9.5.2	Bad Data Channel specification	27
9.5.3	Concatenate not available	27
9.5.4	Source_Concatenate not available	27
9.5.5	Persistent not available	27
9.5.6	Out of Range B_num, Bufx, Offset, S_count, or Sync	28
9.5.7	Request_To_Receive problem	28
9.5.8	Undefined Flag	28

Tables

Table 1 – Response to a rejected Operation	12
Table 2 – Virtual Connection Operations summary between end devices A and B	24
Table 3 – Data transfer and status Operations summary between end devices G and H.....	25
Table 4 – Operation pairs guarded by Op_timeout	26
Table 5 – Summary of logged errors	28

Figures

Figure 1 – System overview.....	4
Figure 2 – HIPPI-ST over different media	4
Figure 3 – Information hierarchy	4
Figure 4 – Scheduled Transfer Final Destination model	5
Figure 5 – Data packing examples	11
Figure 6 – HIPPI-ST service interface.....	13
Figure 7 – Schedule Header contents	14
Figure 8 – Flags summary	15
Figure A.1 – HIPPI-ST Operations carried in HIPPI-6400-PH Messages	30
Figure A.2 – HIPPI-ST Operations carried in HIPPI-FP packets	31
Figure B.1 – Many-to-one striping	33
Figure B.2 – One-to-many striping	33
Figure B.3 – Many-to-many striping	33

Annexes

A Using lower layer protocols	29
A.1 HIPPI-6400-PH as the lower layer.....	29
A.2 HIPPI-FP as the lower layer	29
B HIPPI-ST striping.....	32
B.1 Striping principles	32
B.2 Many-to-one striping	32
B.3 One-to-many striping	32
B.4 Many-to-many striping	33
C Scheduled Transfer example.....	34

Foreword (This foreword is not part of American National Standard X3.xxx-199x.)

This American National Standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

This standard provides an upward growth path for legacy HIPPI-based systems.

This document includes annexes which are informative and are not considered part of the standard.

Requests for interpretation, suggestions for improvement or addenda, or defect reports are welcome. They should be sent to the X3 Secretariat, Information Technology Industry Council, 1250 Eye Street, NW, Suite 200, Washington, DC 20005.

This standard was processed and approved for submittal to ANSI by Accredited Standards Committee on Information Processing Systems, X3. Committee approval of the standard does not necessarily imply that all committee members voted for approval. At the time it approved this standard, the X3 Committee had the following members:

(List of X3 Committee members to be included in the published standard by the ANSI Editor.)

Subcommittee X3T11 on Device Level Interfaces, which developed this standard, had the following participants:

(List of X3T11 Committee members, and other active participants, at the time the document is forwarded for public review, will be included by the Technical Editor.)

Introduction

This American National Standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

Characteristics of a HIPPI-ST include:

- A hierarchy of data units (Scheduled Transfer Units (STUs), Blocks, and Transfers).
- Support for Get and Put Operations.
- Parameters exchanged between end devices for port selection, transfer identification, and Operation validation.
- Features supporting efficient mapping between the sender's and receiver's natural buffer sizes.
- Provisions for resending partial Transfers for error recovery.
- Mappings onto HIPPI-6400-PH, HIPPI-FP (for HIPPI-800 traffic), and Ethernet lower-layer protocols.
- Mappings from IPv4, IPv6, and MPI upper-layer protocols onto Scheduled Transfer.

American National Standard for Information Technology –

High-Performance Parallel Interface – Scheduled Transfer (HIPPI-ST)

1 Scope

This American National Standard specifies a data transfer protocol that uses small control messages to pre-arrange data movement. Buffers are allocated at each end before the data transmission, allowing full-rate, non-congesting data flow between the end devices. The control and data may use different physical media, or may share a single physical medium. Procedures are provided for moving data over HIPPI and other media.

Specifications are included for:

- Virtual Connection **set up** and teardown;
- determining the number of Operations the other end can accept;
- determining the buffer size of the other end;
- exchanging Key, Port, transfer identifiers, and buffer size values specific to the end nodes;
- determining a maximum size **transmission unit** that will not overrun receiver buffer boundaries;
- using buffer indices and 64-bit addresses;
- acknowledging partial transfers so that buffers can be reused;
- providing means for resending partial Transfers for error recovery; and
- terminating transfers in progress.

Note that parts of the Scheduled Transfer protocol depend upon in-order delivery by the lower layer, which may not be available on all media.

2 Normative references

The following American National Standards contain provisions which, through reference in this text, constitute provisions of this American National Standard. At the time of publication, the editions indicated were valid. All standards are subject to revision, and parties to agreements based on this standard are encouraged to investigate the possibility of applying the most recent editions of the standards listed below.

ANSI X3.183-1991, *High-Performance Parallel Interface – Mechanical, Electrical, and Signalling Protocol Specification (HIPPI-PH)*

ANSI X3.210-1992, *High-Performance Parallel Interface – Framing Protocol (HIPPI-FP)*

ANSI X3.222-1993, *High-Performance Parallel Interface – Physical Switch Control (HIPPI-SC)*

ANSI X3.xxx-199x, *High-Performance Parallel Interface – 6400 Mbit/s Physical Layer (HIPPI-6400-PH)*

ANSI/IEEE Std 802-1990, *IEEE Standards for Local and Metropolitan Area Networks: Overview and architecture (formerly known as IEEE Std 802.1A, Project 802: Local and Metropolitan Area Network Standard — Overview and Architecture)*.

ISO/IEC 8802-2:1989 (ANSI/IEEE Std 802.2-1989), *Information Processing Systems – Local Area Networks – Part 2: Logical link control*.

3 Definitions and conventions

3.1 Definitions

For the purposes of this standard, the following definitions apply.

3.1.1 Block: An ordered set of one or more STUs within a Scheduled Transfer. (See figure 3 and 4.4.4.)

3.1.2 Buffer Index (Bufx): A 32-bit parameter identifying the starting address of a data buffer. Bux may be either a pointer to the starting address, or the most significant part of a 64-bit starting address.

3.1.3 Concatenate: An addressing mode using 64-bit addresses rather than buffer indices. (See 4.3.7.)

3.1.4 Control Channel: The logical channel that carries the Control Operations.

3.1.5 Control Operation: A control function consisting of a Schedule Header and an optional 32-byte payload. (See figure 3.)

3.1.6 Data Channel: The logical channel that carries the data payload.

3.1.7 Data Operation: A data movement Operation consisting of a Schedule Header and up to 2 gigabytes of user payload. (See figure 3).

3.1.8 Final Destination: The end device that receives, and operates on, the data payload. This is typically a host computer system, but may also be a non-transparent translator, bridge, or router.

3.1.9 Key: A local identifier used to validate Operations. (See 4.3.3.)

3.1.10 log: The act of making a record of an event for later use.

3.1.11 Operation: A Scheduled Transfer function, i.e., a Control Operation or the data movement specified in an STU.

3.1.12 optional: Characteristics that are not required by HIPPI-ST. However, if any optional characteristic is implemented, it shall be implemented as defined in HIPPI-ST.

3.1.13 Originating Source: The end device that generates the data payload. This is typically a host computer system, but may also be a non-

transparent translator, bridge, or router.

3.1.14 Persistent: A control mode used to retain buffers for multiple Transfers. (See 4.3.9.)

3.1.15 Port: A logical connection within an end device. (See 4.3.2.)

3.1.16 Scheduled Transfer: An information transfer, normally used for bulk data movement and low processing overhead, where the Originating Source and Final Destination prearrange the transfer using the protocol defined in this standard.

3.1.17 Scheduled Transfer Unit (STU): The data payload portion of a Data Operation moved from an Originating Source to a Final Destination. STUs are the basic components of Blocks. (See figure 3 and 4.4.6.)

3.1.18 Slot: A space reserved for a Control Operation, or the Schedule Header portion of an STU, in the end device. (See 4.3.6.)

3.1.19 Transfer: An ordered set of one or more Blocks within a Scheduled Transfer. (See figure 3 and 4.2.)

3.1.20 upper-layer protocol (ULP): The protocol above the service interface. These could be implemented in hardware, software, or they could be distributed between the two.

3.1.21 Virtual Connection: A bi-directional logical connection used for Scheduled Transfers between two end devices. A Virtual Connection contains a logical Control Channel and a logical Data Channel in each direction.

3.2 Editorial conventions

In this standard, certain terms that are proper names of signals or similar terms are printed in uppercase to avoid possible confusion with other uses of the same words (e.g., **END**). Any lowercase uses of these words have the normal technical English meaning.

A number of conditions, sequence parameters, events, states, or similar terms are printed with the first letter of each word in uppercase and the rest lowercase (e.g., Block, Transfer). Any lowercase uses of these words have the normal technical English meaning.

The word *shall* when used in this American National standard, states a mandatory rule or requirement. The word *should* when used in this standard, states a recommendation.

3.2.1 Binary notation

Binary notation is used to represent relatively short fields. For example a two-bit field containing the binary value of 10 is shown in binary format as b'10'.

3.2.2 Hexadecimal notation

Hexadecimal notation is used to represent some fields. For example a two-byte field containing a binary value of b'11000100 00000011' is shown in hexadecimal format as x'C403'.

3.3 Acronyms and other abbreviations

ACK	acknowledge indication
CTS	Clear_To_Send
ENDA	END_ACK
HIPPI	High-Performance Parallel Interface
K	kilo (2^{10} or 1024)
MAC	Media Access Control
MPI	Message Passing Interface
PT	Port_Teardown
PTA	Port_Teardown_ACK
PTC	Port_Teardown_Complete
RQP	Request_Port
RQPR	Request_Port_Response
RS	Request_State
RTR	Request_To_Receive
RTS	Request_To_Send

RTSR	Request_To_Send_Response
SR	State_Response
STU	Scheduled Transfer Unit
ULP	upper-layer protocol

4 System overview

This clause provides an overview of the structure, concepts, and mechanisms used in Scheduled Transfers. Figure 1 gives an example of Scheduled Transfers being used to communicate between device A and device B over some physical media. Annex C describes the steps in a typical Scheduled Transfer. Figure 2 shows HIPPI-ST being used over different media.

4.1 Control Channels and Data Channels

Each Transfer has an Originating Source and Final Destination. Each Originating Source and Final Destination shall have a Control Channel, and one or more Data Channels. The Originating Source sends the payload data, and the Final Destination receives the payload data.

Control Operations shall be exchanged over the Control Channel. Scheduled Transfer Units (STUs), i.e., data payload, shall be exchanged over the Data Channel(s). The information volume on the Data Channel(s) will be probably many times the volume on the Control Channel, hence the available bandwidths should be balanced accordingly. For best performance, the Control Channel should have low latency.

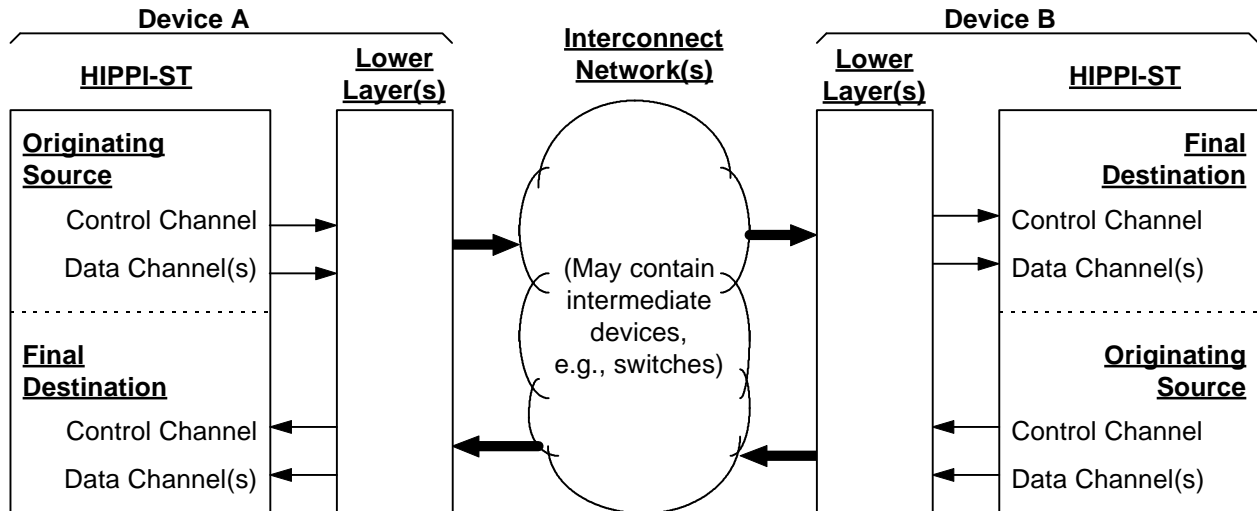


Figure 1 – System overview

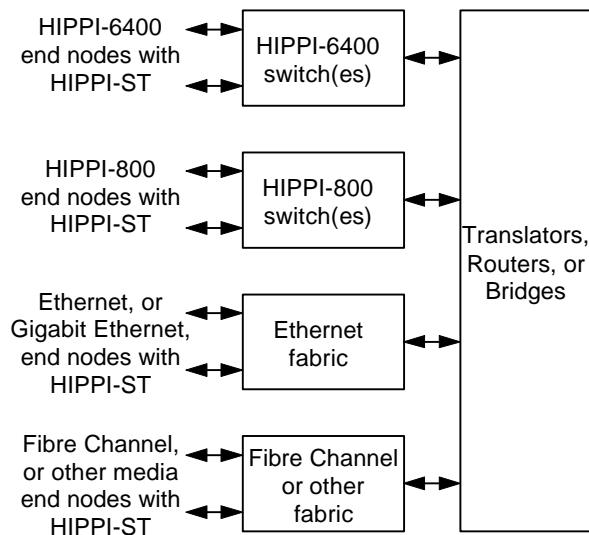
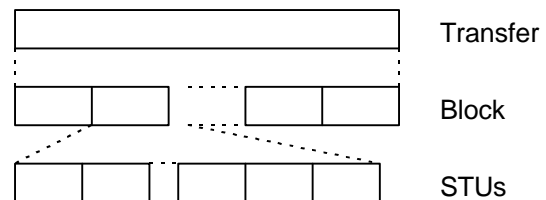


Figure 2 – HIPPI-ST over different media

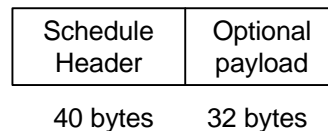
4.2 System model

Multiple write (Put) or read (Get) functions may be executed to move data units, called Transfers, over a Virtual Connection. As shown in figure 3, a Transfer is composed of one or more Blocks, and Blocks are composed of one or more STUs. The Scheduled Transfer protocol shall package the Transfer in Blocks and STUs for delivery using lower layer protocol(s) and media. An STU

shall be the data payload portion of a Data Operation. A Data Operation shall consist of a 40-byte Schedule Header and an STU of up to 2 gigabytes (2^{31} bytes). A Control Operation shall consist of a 40-byte Schedule Header, and may contain an additional 32 bytes of optional payload.



Control Operation



Data Operation

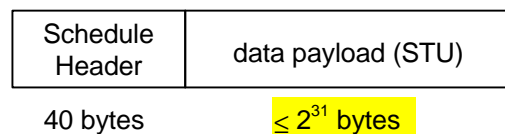


Figure 3 – Information hierarchy

Figure 4 shows the model used on a Final Destination for the Scheduled Transfers. The model on the Originating Source would be similar.

As Control Operations and Data Operations are received, the Schedule Header of each is placed in the Schedule Header queue for execution. State information about the number of empty Slots in the queue is available to the other end so that it can avoid overrunning the queue.

The Virtual Connection Descriptor contains:

- static parameters defining the Virtual Connection from the view of both the remote end device and local end device (the top portion of the Virtual Connection Descriptor box in figure 4);
- current state information about the number of empty "Slots" for Operation Schedule Headers, and Operation Retry **and Timeout** parameters;
- identifiers for each of the Virtual Connection's Transfers.

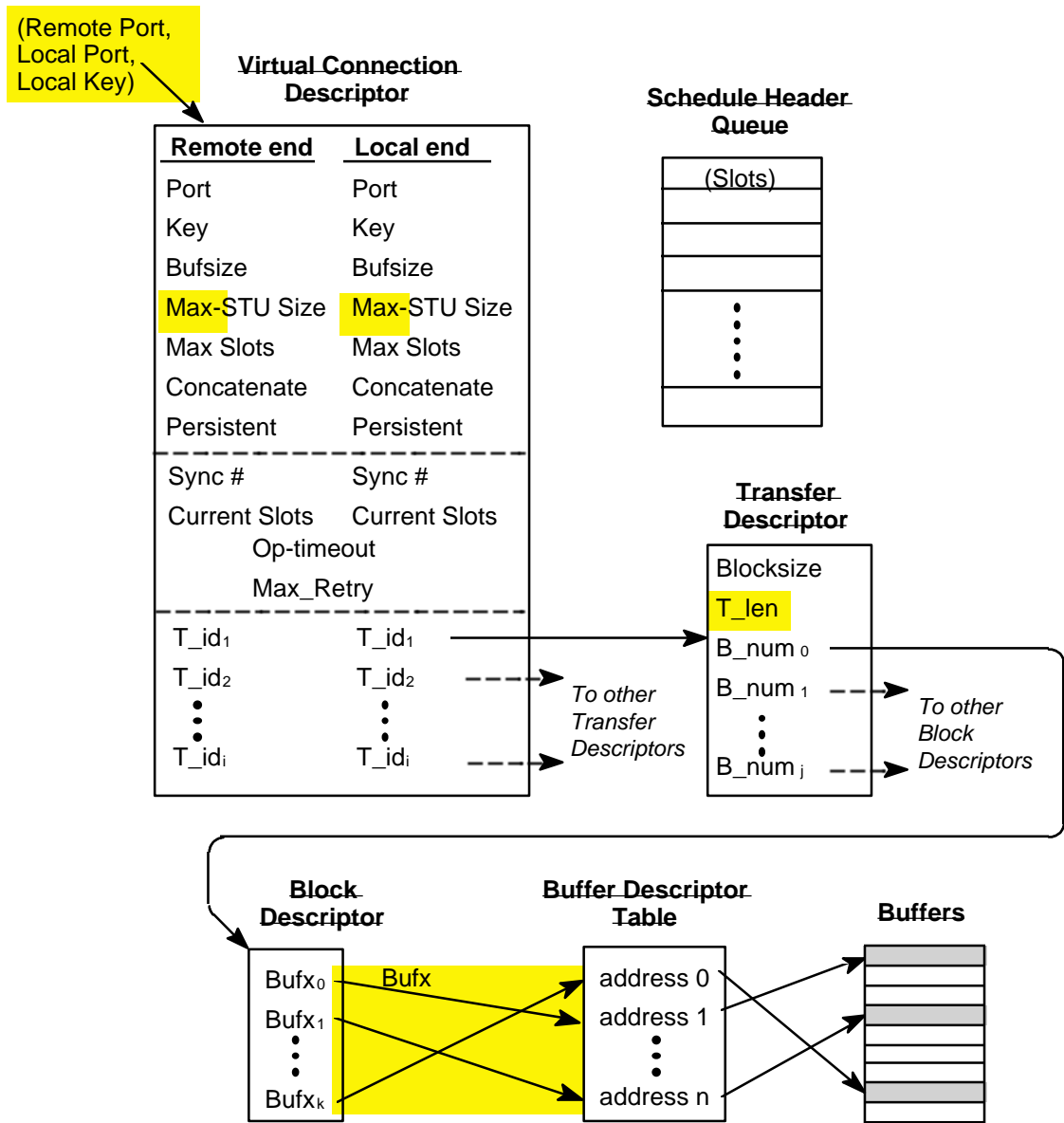


Figure 4 – Scheduled Transfer Final Destination model

A Transfer Descriptor, for each Transfer, contains the Transfer length (T_len, in bytes), the Blocksize (in bytes), and includes pointers to Block Descriptors. The Block Descriptors (one for each Block of a Transfer) identify the set of contiguous Buffer Index (Bufx) values assigned to the Block. And finally, the Buffer Descriptor Table provides a base memory address for each Bufx.

In an effort to achieve maximum transfer rates and efficiency, the receiver's job is made as easy as possible, even at the expense of the transmit side. It is expected that after validating an Operation in the Final Destination, only a single lookup will be needed to derive the absolute memory address and correctly place the data.

4.3 Virtual Connections

Scheduled Transfers between an Originating Source and Final Destination are pre-arranged to decrease computational overhead during the Transfer by allocating buffers at each end device. The bi-directional path between the end devices is called a Virtual Connection. A Virtual Connection shall consist of an Originating Source and Final Destination in each end device.

Once the Final Destination has indicated its ability to accept the STUs, the Virtual Connection should not become congested. In essence, the Final Destination smoothly controls the flow. For comparison, without pre-arranging the buffers, the Originating Source would blindly send data into the interconnection network where it might have to wait for buffers to be assigned in the Final Destination. On the down-side, Scheduled Transfers require additional Control Operations and round-trip latency. Once established, a Virtual Connection may be used to carry multiple Transfers. This Scheduled Transfer protocol does not handle network resource reservations.

4.3.1 Sequences and Operations

During Virtual Connection set up, the end devices shall exchange parameters specific to each device. These parameters, shown in the upper portion of the Virtual Connection Descriptor box in figure 4 and detailed below, include values for:

- Port numbers (e.g., a Port dedicated to HIPPI-FP or IP traffic);
- Keys (used for authenticating Operations);
- native buffer sizes (Bufsize) for determining Final Destination buffer tiling;
- maximum STU size;
- maximum number of outstanding Operations (Slots) to keep from overflowing the command queues;
- whether or not they support Concatenate and Persistent modes.

The parameters assigned during set up shall apply for the life of the Virtual Connection. Once established, the Virtual Connection is accessed as shown in figure 4 by the tuple "remote Port", "local Port", and "local Key". The Control Operations defined for Virtual Connection set up are:

- Request_Port (See 7.1.)
- Request_Port_Response (See 7.2.)

The Control Operations defined for Virtual Connection teardown are:

- Port_Teardown (See 7.3.)
- Port_Teardown_ACK (See 7.4.)
- Port_Teardown_Complete (See 7.5.)

4.3.2 Ports

Ports identify upper-layer entities within an end device. The Port values shall be assigned by the local end device and have no meaning on the other end device. For example, when end device A requests a Virtual Connection to end device B, A shall select the value for A-Port and shall send it to B in the Request_Port Operation. B shall store the A-Port value and shall return it to A in every Operation over this Virtual Connection. Likewise, B shall select the value for B-Port.

An exception is the "well-known Port", i.e., Port x'0000'. In this case, a request sent to the "well-known Port" shall result in the receiving end device assigning a specific local Port value based on the EtherType parameter. EtherType parameter values shall be as assigned in the current "Assigned Numbers" RFC, e.g., RFC 1700¹⁾. For example, if the HIPPI-ST is used to

encapsulate TCP/IP, then the EtherType would be x'0800'. If HIPPI-ST is being used to encapsulate legacy HIPPI-FP or user data, then the EtherTypes would be x'8180' and x'8181' respectively.

If the incoming Port number is invalid, then the Operations shall not be executed (see 9.4.1).

4.3.3 Keys

Like the Ports, each end device shall select its own 32-bit Key value for use on the Virtual Connection. For example, when end device A requests a Virtual Connection to end device B, A shall select the value for A-Key and shall send it to B in the Request_Port Operation. B shall store the A-Key value and shall return it to A in every Operation over this Virtual Connection. The A-Key value has no meaning in B; it is only significant in A where it shall be used to validate that the Operation presented is really associated with this Virtual Connection. Likewise, B shall select the value for B-Key. Keys are similar in nature to passwords; if the Key doesn't match, then the Operation shall not be executed (see 9.4.1).

4.3.4 Buffer size (Bufsize)

Each end shall define the buffer size, in bytes, that it wants to use. Buffer sizes may be the same as host page sizes. It is most efficient when the buffer sizes are the same on both ends, but differing buffer sizes are supported (see annex C). The buffer sizes shall be ≥ 256 bytes and shall be an integral power of two, i.e., 2^x where $8 \leq x \leq 32$.

4.3.5 Max-STU size

The Max-STU size, exchanged during Virtual Connection set up, establishes the maximum data payload size of an STU (see 4.4.9). Each end device declares the desired Max-STU size it is prepared to receive. The Max-STU size must

be no larger than its Bufsize. Intermediate devices with smaller buffer sizes may lower this value. Note that the Max-STU size in each direction may be different.

Additionally, the Max-STU size shall be ≥ 256 bytes and an integral power of two i.e., 2^x where $8 \leq x \leq 32$.

4.3.6 Slots and Sync parameter

The term Slot denotes memory at an end device reserved for storing the Schedule Headers of incoming Operations. Each Operation arriving at an end device consumes one Slot, except for Data Operations which consume a Slot only if a Silent = 0 or Interrupt = 1. An Originating Source shall control the flow of Operations by sending no more Operations than there are Slots available at the other end. Any Operations that are sent in excess of the number of available Slots may be discarded by the receiver (see 9.4.2). The end device supplying the Slot information shall advertise at least one too few Slots, i.e., keeping one in reserve for an END or State_Response Operation to prevent deadlocks.

An end device learns the initial number of Slots available (Slots value) at the remote end device during the Virtual Connection set up (see 7.1 and 7.2). Later, an end device obtains the current Slots value by reading the Slots parameter in a received State_Response. An end device may solicit a State_Response from the remote end by either of two methods: by setting the Send_State flag in the Schedule Header of a Data Operation, or by sending a Request_State Operation. A received Slots value of x'FFFFFFFF' indicates that the remote end does not implement Slot accounting.

NOTE – Slot flow control may not be needed when the maximum number of Control Operations is otherwise bounded or where dropped Operations are acceptable.

¹⁾ RFC (Request For Comment) documents are working standards documents from the TCP/IP internetworking community. Copies of these documents are available from numerous electronic sources (e.g., <http://www.ietf.org>) or by writing to IETF Secretariat, c/o Corporation for National Research Initiatives, 1895 Preston White Drive, Suite 100 Reston, VA 20191-5434, USA.

The received Slot value is a snapshot of the number of Slots available at the remote end device when the remote end device received the soliciting Operation. The local end device may continue to send Operations after soliciting a State_Response and may also solicit multiple responses before receiving a reply. The lower bound on the number of available Slots at the remote end device is determined by the local end device which adjusts its vision of the number of Slots to account for outstanding Operations. The adjustment consists of subtracting, from the number of Slots indicated in the received State_Response Operation, the number of Slot-consuming Operations sent by the local end device after a State_Response solicitation.

The local end device can use the Sync parameter to identify State_Response messages when there are multiple outstanding solicitations. The Sync parameter in a Data or Request_State shall be copied and returned by the remote end device in the corresponding State_Response. The Sync parameter may be used by the local end device to mark the request, and thus identify the State_Response with a particular solicitation. The Sync values are locally determined.

4.3.7 Concatenate

The Concatenate flag (see 6.2) controls the addressing mode for the Bufx and Offset parameters.

- When Concatenate = 0, Bufx shall specify a Buffer Index for placing the data in the Final Destination.
- When Concatenate = 1, the Bufx and Offset fields shall be concatenated into a single 64-bit address, with Bufx containing the most-significant bytes of the address (see 4.4.7).

The value of the Concatenate flag shall be consistent for an entire Transfer, i.e., switching back and forth between 64-bit addresses and Buffer Indexes within a Transfer is not allowed (see 4.4.7).

Concatenate is only usable between hosts that mutually agree. Agreement is reached by controlling the Concatenate flag bit during the Virtual Connection set up (see 7.1 and 7.2).

4.3.8 Source_Concatenate

The Source_Concatenate flag (see 6.2) controls the addressing mode for the OS_Bufx (Originating Source Buffer index) and OS_Offset parameters. Note that OS_Bufx and OS_Offset control the data addressing on the Originating Source and are only used in the Request_To_Receive Operation.

- When Source_Concatenate = 0, OS_Bufx shall specify a Buffer Index pointing to the data in the Originating Source (see 4.4.8).
- When Source_Concatenate = 1, the OS_Bufx and OS_Offset fields shall be concatenated into a single 64-bit address, with OS_Bufx containing the most-significant bytes of the address (see 4.4.8).

Source_Concatenate is only usable between hosts that mutually agree. Agreement is reached by controlling the Source_Concatenate flag bit during the Virtual Connection set up (see 7.1 and 7.2).

4.3.9 Persistent

The Persistent flag (see 6.2) controls buffer retention in the Final Destination for the Virtual Connection.

- When Persistent = 1, the memory in the Final Destination allocated for the Scheduled Transfer shall be retained for multiple transfers and not released until a Port_Teardown or an END Operation occurs. Note that Persistent = 1 bypasses the flow control provided by Clear_To_Send, i.e., a Data Operation may be sent at any time whether or not a Clear_To_Send Operation has been received. Sending information to a Frame Buffer is an example of where Persistent might be used.
- When Persistent = 0, the memory for a Block may be allocated for other uses after the Block is complete. All Data Operations must be enabled by a Clear_To_Send or Request_To_Receive Operation.

Persistent is only usable between hosts that mutually agree. Agreement is reached by controlling the Persistent flag bit during the Virtual Connection set up (see 7.1 and 7.2).

4.4 Data movement

4.4.1 Sequences and Operations

A write data sequence (which may be initiated by either end of the Virtual Connection) shall be set up by the end devices exchanging transfer identifiers (T_id's), specific to each device, and length parameters. The Control Operations setting up a write data sequence are:

- Request_To_Send (See 8.1.)
- RTS_Response (See 8.2.)

A read data sequence, which moves the Transfer as a single Block, requires that both ends had previously allocated resources for the entire read sequence with a Request_To_Send. The Control Operation setting up a read data sequence is:

- Request_To_Receive (See 8.3.)

The Final Destination controls the data flow with:

- Clear_To_Send (See 8.4.)

Data payloads for the read and write data movements are carried in STUs. STUs are sent with:

- Data (See 8.5.)

State information can be requested in a Data Operation or with a Request_State Control Operation.

- Request_State (See 8.6.)
- State_Response (See 8.7.)

The Control Operations below are used to abort limited size Transfers. Unlimited size Transfers shall use this method to signal the end of the Transfer.

- END (See 8.8.)
- END_ACK (See 8.9.)

4.4.2 Transfer identifiers (R_id and S_id)

Like the Ports and Keys, each end device shall also select its own non-zero 16-bit Transfer identifier (T_id) value for a data movement on the Virtual Connection. For example, when end device S requests to write to end device R, S shall select the value for its T_id and shall send it

to R in the Request_To_Send Operation. R shall store S's T_id value and shall return it to S in every Operation concerning this Transfer. Likewise, R shall select its T_id value and send it to S in a Request_To_Send_Response or Clear_To_Send Operation. For each Operation, the sender shall put its T_id in the S_id field and the receiver's T_id value in the R_id field.

NOTE – The Virtual Connection is symmetrical; either end device may initiate a data movement. For example, S could be end device A that initiated the Virtual Connection set up, or it could be end device B. Different names were used for clarity.

4.4.3 Transfer length (T_len)

The 32-bit Transfer length parameter (T_len) specifies the total number of data payload bytes in the Transfer. T_len does not include the Schedule Header or any lower-layer headers. T_len = x'00000000' shall indicate an unlimited size Transfer. An unlimited size Transfer is terminated by an END Operation (see 8.8).

4.4.4 Blocks

Scheduled Transfer flow control, striping, acknowledgments, and resource allocation are all done on a Block basis. Block numbers (B_num) shall be numbered starting at zero and shall increment by one for each following Block.

Blocks comprising a Transfer shall be enabled for transmission in sequential order unless both the Originating Source and Final Destination indicated Out_of_Order capability during the Virtual Connection set up. Note that Out_of_Order is necessary for retransmission to correct flawed Blocks.

State_Response Operations indicate the highest numbered Block received correctly by the Final Destination. State_Response Operations can be requested by setting the Send_State flag bit in Data Operations or by sending Request_State Operations. In addition, Request_State Operations can ask if a particular Block was received correctly. Use of these mechanisms allows the Originating Source to verify correct reception and to identify flawed Blocks for potential retransmission.

4.4.5 Blocksize

The **Blocksize** (in bytes) for a Transfer is established when a Transfer is initiated, i.e., with a **Request_To_Send_Response** or **Clear_To_Send** Operation (see 8.2 and 8.4). **Blocksize** is expressed as a power of two, i.e., 2^x where $8 \leq x \leq 32$. All of the Blocks of a Transfer shall be full size, except for the first and/or last Block of a Transfer which can be smaller (the first Block will be smaller by the initial Offset value, and the last Block will be whatever completes the Transfer).

4.4.6 STUs

The STUs of a Block shall be transmitted in order. STU numbers (**S_count**) shall start with **zero** and **increment** by one for each following STU. The **last** STU of a Block shall **be marked with Last = 1**. No STU shall extend past a Final Destination's **buffer boundary**, **Blocksize** boundary, or Transfer boundary.

4.4.7 Bufx and Offset

Bufx contains either a Buffer Index, or the high-order portion of a 64-bit address, in the Final Destination. Selection between these two is controlled by the Concatenate flag (see 4.3.7 and 6.2). If more than one Buffer Index is required for a Block, i.e., buffer size (Bufsize) is less than **Blocksize**, then the Bufx parameter in the **Clear_To_Send** Operation shall specify the initial Bufx, and any additional Bufx values shall be sequential.

Offset may be used to start at other than the first byte of a Final Destination's buffer. For the first STU of a Block, the Offset value shall be the same as received in the **Clear_To_Send** for the Block. Subsequent STUs of the Block shall adjust the Bufx and Offset values based on the Final Destination's buffer size and the STU size used by the Originating Source.

The Offset value associated with the first block of a Transfer (**I_Offset**) is included in all **Clear_To_Send** Operations. This allows the Originating Source to compute the starting address for any Block without having received the **Clear_To_Send** for the first Block.

Clear_To_Send Operations can occur out of order, e.g., as the result of striping.

Best performance will usually be achieved when an Offset value of zero is specified. Use of non-zero offset values may degrade performance, depending upon underlying hardware transfer mechanisms.

4.4.8 OS_Bufx and OS_Offset

OS_Bufx specifies either a Buffer Index, or the high-order portion of a 64-bit address, in the Originating Source during a **Request_To_Receive** Operation. Selection between these two is controlled by the Source_Concatenate flag (see 4.3.8 and 6.2). If more than one Buffer Index is required for a Block, i.e., buffer size < **Blocksize**, then the OS_Bufx parameter shall specify the initial Bufx, and any additional Bufx values shall be sequential.

OS_Offset may be used to start at other than the first byte of a Source buffer. Note that OS_Bufx and OS_Offset are only used with **Request_To_Receive** Operations, and **Request_To_Receive** Operations only specify one Block.

4.4.9 Packing examples

Figure 5 shows three possibilities for packing the same Transfer into a receiver's buffers. All three examples show a group of seven of the receiver's buffers on the top line. Each buffer is pointed to by a Bufx, and the data in the first buffer starts at an Offset value. The Transfer is the shaded bar, with transmission going from left to right. The Block boundaries are shown above the shaded bar, and the resulting STU boundaries are shown below the shaded bar.

Example (a), at the top, shows the case where the buffers and Blocks are the same size. Notice that the first Block is smaller than the other Blocks by the Offset value. No Offset is required for the other Blocks. The last Block of the Transfer is also smaller, i.e., the Transfer did not end on a Block boundary. While the STU boundaries lined up nicely, the sender could have used multiple STUs, but the STUs cannot be larger than Max-STU.

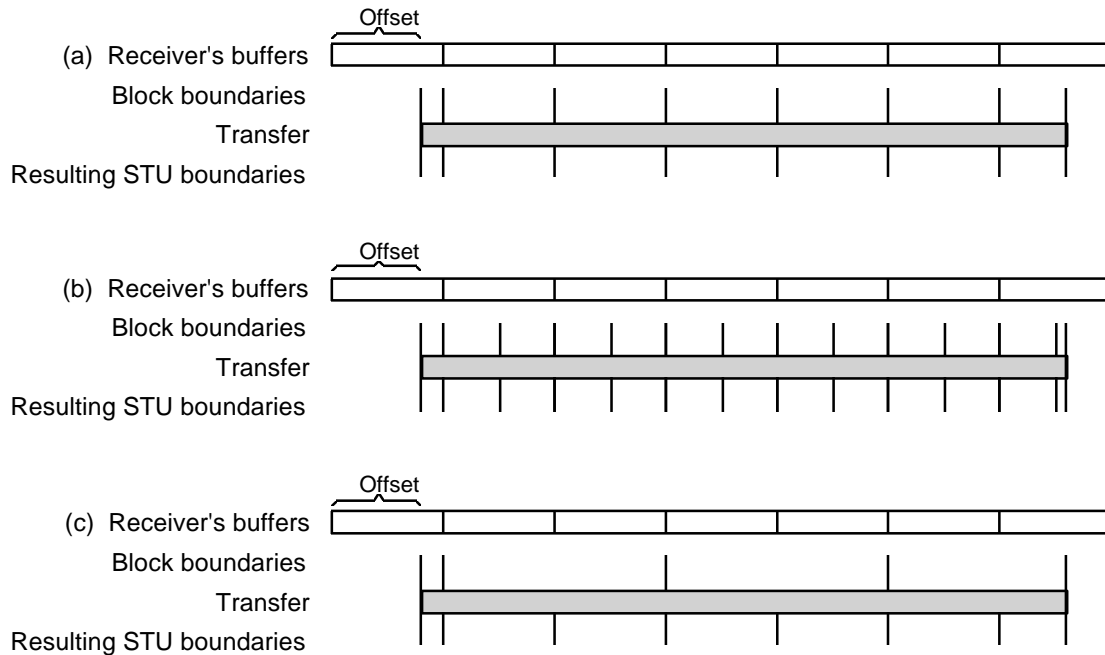


Figure 5 – Data packing examples

Example (b) shows multiple Blocks per receiver buffer. The Blocks that do not start on a buffer boundary would use the Offset parameter to position the data.

Example (c) shows the Blocksize covering two of the receiver's buffers.

In summary, STUs cannot cross Block, buffer, or Transfer boundaries. Relationships include:

STU size \leq Max-STU size

STU size \leq Blocksize

STU size \leq Bufsize

Max-STU size \leq Blocksize

Max-STU size \leq Bufsize

Note that Blocksize can be larger, smaller, or the same as Bufsize.

4.5 Operations management

4.5.1 Flow control

Data flow control is achieved with Clear_To_Send and Request_To_Receive Operations; each one sent by the Final Destination gives the Originating Source permission to send one Block. Flow control is over-ridden when Persistent = 1; here Data Operations may be sent without having first received Clear_To_Send Operations.

Operation flow control is achieved by an Operation's sender not overrunning the Slots value (see 4.3.6).

4.5.2 Status Operations

Request_State (see 8.6) and State_Response (see 8.7) Operations are used to request and supply status information about the state of the remote end device. They can be used to see which Blocks have been received correctly and the number of empty Slots available. The Sync parameter (see 4.3.6) is used to provide a common reference point for the local and remote end devices, i.e., to match State_Request and State_Response Operations.

4.5.3 Rejected Operations

If the receiving end device is unable to execute an Operation, then the receiving device shall set the Reject flag bit = 1 in the response. Table 1 shows the response when an Operation is rejected. The recovery actions taken when an Operation is rejected are beyond the scope of this standard.

Table 1 – Response to a rejected Operation

Rejected Operation	Response (w/ Reject = 1)
Request_Port	Request_Port_Response
Request_To_Send	Request_To_Send_Response
Request_To_Receive	State_Response

4.5.4 Lost Operations

Errors other than syntactic errors are manifested as missing Operations, which occur when the underlying physical medium discards or damages a transmission. Each Scheduled Transfer Operation is defined as part of a two-way handshake or a three-way handshake. Thus, for each command Operation there is a corresponding response Operation, and for some response Operations there is also a corresponding completion Operation.

Each Operation that expects a response is guarded with a timeout whose value is referred to as Op_timeout (see 9.1). An Operation shall be re-tried up to Max_Retry times (see 9.1) if the sending end device does not receive the expected response (see 9.2 and table 5).

Data transmissions (i.e., Data Operations) are an exception to this timeout mechanism and are referred to the ULP for resolution (see 9.2).

4.5.5 Interrupts

An Interrupt causes a signal to be delivered to the receiving end device ULP. An Interrupt can be requested with any Operation by setting Interrupt = 1.

5 Service interface

This clause specifies the services provided by HIPPI-ST. The intent is to allow ULPs to operate correctly with this HIPPI-ST. How many of the services described herein are chosen for a given implementation is up to that implementor; however, a set of HIPPI-ST services must be supplied sufficient to satisfy the ULP(s) being used. The services as defined herein do not imply any particular implementation or any interface.

Figure 6 shows the relationship of the HIPPI-ST interfaces.

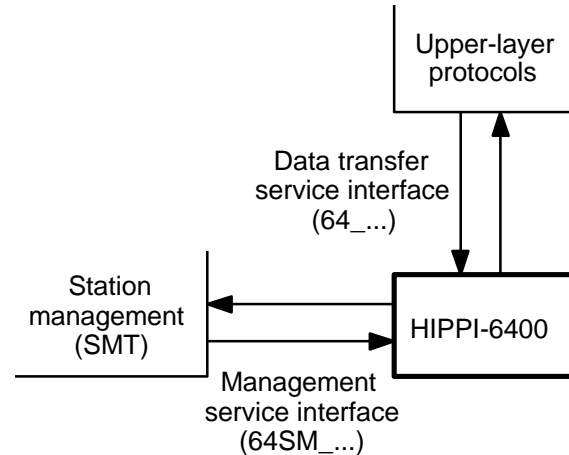


Figure 6 – HIPPI-ST service interface

5.1 Service primitives

The primitives, in the context of the state transitions in clause 5, are declared required or optional. Additionally, parameters are either required, conditional, or optional. All of the primitives and parameters are considered as required except where explicitly stated otherwise.

HIPPI-ST service primitives are of four types.

- *Request primitives* are issued by a service user to initiate a service provided by the HIPPI-ST. In this standard, a second Request primitive of the same name shall not be issued until the Confirm for the first request is received.
- *Confirm primitives* are issued by the HIPPI-ST to acknowledge a Request.
- *Indicate primitives* are issued by the HIPPI-ST to notify the service user of a local event. This primitive is similar in nature to an unsolicited interrupt. Note that the local event may have been caused by a service Request. In this standard, a second Indicate primitive of the same name shall not be issued until the Response for the first Indicate is received.
- *Response primitives* are issued by a service user to acknowledge an Indicate.

5.2 Sequences of primitives

The order of execution of service primitives is not arbitrary. Logical and time sequence relationships exist for all described service primitives. Time sequence diagrams are used to illustrate a valid sequence. Other valid sequences may exist. The sequence of events between peer users across the user/provider interface is illustrated. In the time sequence diagrams, the HIPPI-ST users are depicted on either side of the vertical bars, while the HIPPI-ST acts as the service provider.

NOTE - The intent is to flesh out the service primitives similar to what is in HIPPI-PH today.

Service interface considerations -

(These are notes that have been collected during the document reviews, and should be considered when the service interface is written.)

Should there be a priority, or time-to-live, for individual Transfers? On a per connection basis?

Pass the full ST header to/from the ULP.

Service the slots in order of arrival, i.e., FIFO.

Interrupts are passed independent of the Slots, i.e., whenever an Interrupt is put in the Slots queue.

There is a Port-basis for the Service Interface

6 Schedule Header

The Schedule Header is shown in figure 7 as a group of 32-bit words. The Schedule Header fields are named for the most common parameter for which the field is used. Many of the fields have different uses depending on the Operation type, and some Operations do not use one or more of the fields at all. The usage for each field is listed below and summarized in tables 2 and 3.

			Bytes
Op	Flags	S_count	00-03
R_Port		S_Port	04-07
Key			08-11
R_id		S_id	12-15
Bufx			16-19
Offset			20-23
T_len			24-27
B_num			28-31
OS_Bufx			32-35
OS_Offset			36-39

Figure 7 – Schedule Header contents

6.1 Schedule Header fields

The Schedule Header fields shall be as follows. If an Operation does not use a particular Schedule Header field, then that field shall be transmitted as zeros.

Op (5 bits, high-order 5 bits of byte 00) – The Scheduled Transfer Operation. See tables 2 and 3 for a summary of Op values. Unspecified Op values are reserved.

Flags (11 bits, low-order 3 bits of byte 00, and all of byte 01) – Control flags (see 6.2).

S_count (16 bits, bytes 02-03):

- In *Request_Port*, *Request_Port_Response*, and *State_Response* Operations: the number of available Slots (see 4.3.6);
- In *Request_To_Send_Response* and *Clear_To_Send* Operations: the Blocksize parameter (see 4.4.5);

– In *Data* Operations: the STU number (see 4.4.6).

R_Port (16 bits, bytes 04-05) – The receiver's logical Port for this Operation (see 4.3.2).

S_Port (16 bits, bytes 06-07) – The sender's logical Port for this Operation (see 4.3.2).

Key (32 bits, bytes 08-11) – Virtual Connection identifier. Generated independently by each end during the Virtual Connection set up. (See 4.3.3.)

R_id (16 bits, bytes 12-13) – The receiver's Transfer identifier for this Operation (see 4.4.2).

S_id (16 bits, bytes 14-15) – The sender's Transfer identifier for this Operation (see 4.4.2).

Bufx (32 bits, bytes 16-19):

– In *Request_Port* and *Request_Port_Response* Operations: the maximum buffer size (Bufsize) supported by the end device (see 4.3.4);

– In *Request_To_Receive*, *Clear_To_Send*, and *Data* Operations: the Buffer Index at the Final Destination or the high-order portion of a 64-bit concatenated address (see 4.4.7).

Offset (32 bits, bytes 20-23):

– In *Request_Port* and *Request_Port_Response* Operations: the sender's Key value (see 4.3.3);

– In *Request_To_Receive*, *Clear_To_Send*, and *Data* Operations: the Final Destination's Offset within a Bufox, or the low-order portion of a 64-bit concatenated address (see 4.4.7);

– In *State_Response* Operations: the Block number of the highest numbered contiguous Block received correctly (see 4.4.4).

T_len (32 bits, bytes 24-27):

– In *Request_Port* and *Request_Port_Response* Operations: the Max-STU size (see 4.3.5);

– In *Request_To_Send*, *Request_To_Send_Response*, and *Request_To_Receive* Operations: the length, in bytes, of the Transfer data (see 4.4.3);

– In *Data*, *State_Request*, and *State_Response* Operations: the Sync parameter (see 4.3.6).

B_num (32 bits, bytes 28-31):

- In *Request_Port Operations*: the EtherType parameter (see 4.3.2);
- In *Clear_To_Send* and *Data Operations*: the Block number being requested or transmitted (see 4.4.4);
- In *Request_State* and *State_Response Operations*: the Block number being queried or responded to (see 4.4.4, 8.6, and 8.7).

OS_Bufx (32 bits, bytes 32-35):

- In *Request_To_Receive Operations*: the Buffer Index at the Originating Source or the high-order portion of a 64-bit concatenated address (see 4.4.8);
- In *Data Operations*: Opaque field, e.g., for passing ULP parameters.

OS_Offset (32 bits, bytes 36-39):

- In *Request_To_Receive Operations*: the Originating Source's Offset within a Bufx, or the low-order portion of a 64-bit concatenated address (see 4.4.8);
- In *Clear_To_Send Operations*: the Final Destination's initial Offset value (see 4.4.7);
- In *Data Operations*: Opaque field, e.g., for passing ULP parameters.

6.2 Scheduled Transfer flags

Figure 8 summarizes the flags, and shows their relative position. The flag functions are detailed below for the case where the bit = 1.

Out_of_Order (b'1xxxxxxxxx') = The end device is able to send and receive Blocks in any order.

Silent (b'x1xxxxxxxxx') = Requests silent delivery of a Data Operation. For Control Operations the Silent flag shall be ignored (i.e., transmitted as zero, but not checked at the receiver). For Data Operations with Silent = 1 the data transfer to the Destination OS_Bufx is carried out normally, but the Schedule Header shall not be delivered to any upper-layer entity. This provides the basis for remote memory write semantics where the intent is to modify the contents of a remote memory without executing software in the Destination host computer.

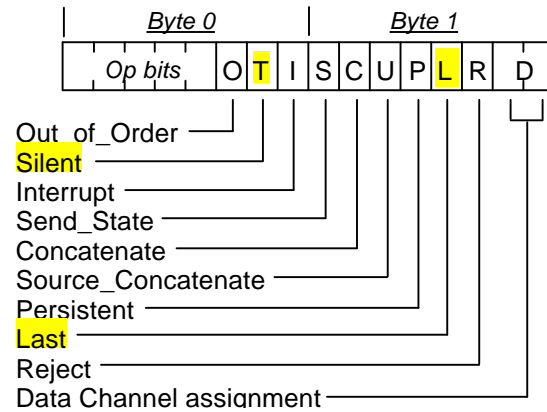


Figure 8 – Flags summary

Interrupt (b'xx1xxxxxxxx') = Requests that a signal or interrupt be generated and delivered to the appropriate upper-layer entity. The Interrupt flag is independent of the Silent flag, i.e., Interrupt = 1 calls for a signal whether or not Silent = 1. (See 4.5.5.)

NOTE 1 – The Silent and Interrupt flags together provide for three delivery modes for Data Operations: silent, polled, or interrupt-driven. If Silent = 1, the data payloads are delivered silently. If Silent = 0, then the upper-layer entity is informed by the same means used for all other Schedule Headers. This mode is suitable for polled interfaces. If Interrupt = 1, then a signal is delivered.

Send_State (b'xxx1xxxxxxxx') = Requests that the Final Destination respond with a State_Response upon successful receipt of this STU by the higher-layer protocol. For Send_State to be valid, either Interrupt = 1 or Silent = 0 must be true.

Concatenate (b'xxxx1xxxxx') = Use Bufx as the high-order portion of a 64-bit Final Destination address rather than as a Buffer Index (see 4.3.7).

Source_Concatenate (b'xxxxx1xxxxx') = Use OS_Bufx as the high-order portion of a 64-bit Originating Source address rather than as a Buffer Index (see 4.3.8).

Persistent (b'xxxxxx1xxxx') = Retain the Final Destination's buffers (see 4.3.9).

Last (b'xxxxxxx1xxx') = The last STU of a Block.

Reject (b'xxxxxxxx1xx') = The request (i.e., Request_Port, Request_To_Send, or Request_To_Receive) has been rejected.

Data Channel assignment: The Data Channel to be used to carry Data Operations. The Data Channel value is assigned in a Request_To_Send Operation and is the Data Channel to be used for Data Operations associated with this Transfer.

b'xxxxxxxx01' = Data Channel 1

b'xxxxxxxx10' = Data Channel 2

b'xxxxxxxx11' = Data Channel 3

The maximum STU size sent on Data Channels 1 and 2 shall be 2^{17} bytes (i.e., 128 Kbytes). The maximum STU size sent on Data Channel 3 shall be 2^{31} bytes (i.e., 2 gigabytes).

NOTE 2 – Data Channel assignment value b'00' is reserved.

7 Virtual Connection management

In this clause, a Virtual Connection is set up between two Ports (see 4.3.2), called the A-Port and B-Port. The device that initiates the Virtual Connection is called device A, and the device at the other end is called device B.

In addition to the Port values, each Port shall assign and associate a Key value (A-Key and B-Key) with the Virtual Connection (see 4.3.3). Other parameters exchanged during the Virtual Connection set up include Buffer sizes (A-BuFSIZE and B-BuFSIZE, see 4.3.4), maximum STU sizes (Max-STU, see 4.3.5), and the number of available Slots (A-Slots and B-Slots, see 4.3.6). The end devices also inform each other of their capability to support Concatenate (see 4.3.7), Source_Concatenate (see 4.3.8), Persistent (see 4.3.9), and out-of-order Block delivery (see 4.4.4).

The Operations used to set up and tear down Virtual Connections are detailed below and summarized in table 2. Only the fields used in each Operation are listed; all of the other Schedule Header fields shall be transmitted as zeros. While a particular field usually carries the parameter of the same name, fields sometimes carry other parameter values. In the Operations below, the specific parameter used in the

Operation is listed first, and if it is not carried in the field of the same name, then the field name is included in square brackets.

7.1 Request_Port

Request_Port shall be used to set up a Virtual Connection between end device A and end device B.

Semantics – Request_Port (

Op,
Flags,
A-Slots [S_count],
B-Port [R_Port],
A-Port [S_Port],
A-BuFSIZE [Bufx],
A-Key [Offset],
A-Max-STU [T_len],
EtherType [B_num])

Op = x'01'

Flags (see 6.2) shall specify the Out_of_Order, Source_Concatenate, Concatenate, and Persistent flags. A value of 1 shall indicate that A supports that feature. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

A-Slots, carried in the S_count field, shall specify the maximum number of Slots allocated in A for this Virtual Connection (see 4.3.6).

B-Port, carried in the R_Port field, shall specify B's logical Port value for this Virtual Connection. B-Port may be either the well-known Port (B will assign the Port value), or a peer Port, that provides the service (see 4.3.2).

A-Port, carried in the S_Port field, shall specify A's logical Port value for this Virtual Connection (see 4.3.2).

A-BuFSIZE, carried in the BuFX field, shall specify A's buffer size (see 4.3.4).

A-Key, carried in the Offset field, shall specify A's Key value for this Virtual Connection (see 4.3.3).

A-Max-STU, carried in the T_len field, shall be \leq A-BuFSIZE when sent by A (see 4.3.5). The A-Max-STU value received by B shall be used by B as the maximum size of STUs (Max-STU) to be sent from B to A on this Virtual Connection. (See 4.3.5.)

EtherType, carried in the B_num field, shall be a value that characterizes the **ULP** data payloads that will be exchanged on this Virtual Connection (see 4.3.2).

Issued – By device A.

Effect – If it accepts the request, then end device B shall establish a Virtual Connection and shall reply with a Request_Port_Response Operation. If rejected, then end device B shall respond with Reject = 1 in the Request_Port_Response (see 4.5.3).

7.2 Request_Port_Response

Request_Port_Response shall inform end device A whether the Virtual Connection was accepted or not. If accepted, the parameters associated with this Virtual Connection are passed to A.

Semantics – Request_Port_Response (

- Op,
- Flags,
- B-Slots [S_count],
- A-Port [R_Port],
- B-Port [S_Port],
- A-Key [Key],
- B-Bufsize [Bufx],
- B-Key [Offset],
- B-Max-STU [T_len])

Op = x'02'

Flags (see 6.2) shall specify the Out_of_Order, Source_Concatenate, Concatenate, and Persistent flags. A value of 1 shall indicate that B supports that feature. The appropriate value for the Reject and Interrupt flags shall also be carried (see 4.5.3 and 4.5.5).

B-Slots, carried in the S_count field, shall specify the maximum number of Slots allocated in B for this Virtual Connection (see 4.3.6).

A-Port, carried in the R_Port field, shall be the same as the A-Port value in the Request_Port Operation (see 4.3.2).

B-Port, carried in the S_Port field, shall specify B's logical Port value for this Virtual Connection (see 4.3.2).

A-Key, carried in the Key field, shall be the Key value assigned by A in the Request_Port Operation (see 4.3.3).

B-Bufsize, carried in the Bufx field, shall specify B's buffer size (see 4.3.4).

B-Key, carried in the Offset field, shall specify B's Key value assigned for this Virtual Connection (see 4.3.3).

B-Max-STU, carried in the T_len field, shall be \leq B-Bufsize when sent by B (see 4.3.5). The **B-Max-STU** value received by A shall be used by A as the maximum size of STUs (Max-STU) to be sent from A to B on this Virtual Connection. (See 4.3.5.)

Issued – By B in response to a Request_Port.

Effect – End device A has been assigned a logical Port on end device B. The Ports, Keys, buffer sizes, maximum STU size, and maximum number of Slots have been exchanged, and a Virtual Connection has been established. Note that the Virtual Connection is bi-directional in that either A or B may initiate a Scheduled Transfer. Multiple Scheduled Transfers may occur over a single Virtual Connection, and the Scheduled Transfers can be either writes or reads.

7.3 Port_Teardown

Port_Teardown shall terminate the Virtual Connection and may be issued by either end device A or end device B. The Port_Teardown sequence uses a three-way handshake consisting of Port_Teardown, Port_Teardown_ACK, and Port_Teardown_Complete.

Open Issue – A state table describing the 3-way handshake will be included in a Normative annex.

Semantics – Port_Teardown (

- Op,
- Flags,
- R_Port,
- S_Port,
- Key)

Op = x'03'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R_Port shall contain the value associated with the receiver of the Operation, e.g., R_Port = B-Port when the Port_Teardown is issued by A (see 4.3.2).

S_Port shall contain the value associated with the sender of the Operation, e.g., S_Port = A-Port when the Port_Teardown is issued by A (see 4.3.2).

Key shall contain the Key value associated with the receiver of the Operation, e.g., Key = B-Key when the Port_Teardown is issued by A (see 4.3.3).

Issued – By either side, i.e., end device A or end device B, of the Virtual Connection. The sender should only issue a Port_Teardown when the Transfers are complete or appear to be stalled.

Effect – The receiver should release any buffers associated with this Virtual Connection, but shall retain the Port and Key values for use in further Port_Teardown Operations. The receiver shall also respond with a Port_Teardown_ACK.

7.4 Port_Teardown_ACK

Port_Teardown_ACK shall be used to acknowledge receipt of a Port_Teardown.

Semantics – Port_Teardown_ACK (

- Op,
- Flags,
- R_Port,
- S_Port,
- Key)

Op = x'04'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R_Port shall contain the value associated with the receiver of the Operation, e.g., R_Port = B-Port when the Port_Teardown_ACK is issued by A (see 4.3.2).

S_Port shall contain the value associated with the sender of the Operation, e.g., S_Port = A-Port when the Port_Teardown_ACK is issued by A (see 4.3.2).

Key shall contain the Key value associated with the receiver of the Operation, e.g., Key = B-Key when the Port_Teardown_ACK is issued by A (see 4.3.3).

Issued – By the receiver of a Port_Teardown Operation after releasing this Virtual Connection's buffers.

Effect – The receiver should release any buffers associated with this Virtual Connection, but shall retain the Port and Key values for use in further Port_Teardown Operations. The receiver shall also respond with a Port_Teardown_Complete.

7.5 Port_Teardown_Complete

Port_Teardown_Complete shall be used to complete a three-way handshake, acknowledging that the actions associated with a Port_Teardown have been completed.

Semantics – Port_Teardown_Complete (

- Op,
- Flags,
- R_Port,
- S_Port,
- Key)

Op = x'05'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R_Port shall contain the value associated with the receiver of the Operation, e.g., R_Port = B-Port when the Port_Teardown_Complete is issued by A (see 4.3.2).

S_Port shall contain the value associated with the sender of the Operation, e.g., S_Port = A-Port when the Port_Teardown_Complete is issued by A (see 4.3.2).

Key shall contain the Key value associated with the receiver of the Operation, e.g., Key = B-Key when the Port_Teardown_Complete is issued by A (see 4.3.3).

Issued – By the receiver of a Port_Teardown_ACK Operation.

Effect – After the **Op_timeout** expires twice, both the sender and receiver shall release the Virtual Connection's Port and Key values. The delay allows for lost or damaged Port_Teardown Operations to be re-issued.

8 Data movement

The Operations used for Scheduled Transfers are detailed below and summarized in table 3. All of the Scheduled Transfer data transfer Operations use the R_Port, S_Port, A-Key, and B-Key values that were assigned during the Virtual Connection **set up** (see 7.1 and 7.2). When end device A issues the Operation:

R_Port = B-Port
S_Port = A-Port
Key = B-Key

Likewise, when end device B issues the Operation:

R_Port = A-Port
S_Port = B-Port
Key = A-Key

For clarity and brevity, these values are not discussed in the individual Operations. All other Schedule Header **fields** that are not listed in a specific Operation shall be transmitted as zeros. While a particular field usually carries the parameter of the same name, fields sometimes carry other parameter values. In the Operations below, the specific parameter used in the Operation is listed first, and if **it** is not carried in the field of the same name, then the field name is included in square brackets.

8.1 Request_To_Send

Request_To_Send asks that space be allocated, **and authorization be given**, for a Transfer. Request_To_Send is issued by the Originating Source to specify the number of data bytes to be sent from the Originating Source to the Final Destination. In addition, the Originating Source shall specify whether 64-bit address or Buffer Indexes are used, whether the Final Destination's buffer should be persistent or discarded after a Block, and the Data Channel assignment for the data transfer. Note that the end device on either end of the Virtual Connection may issue a Request_To_Send. A Request_To_Send, with Persistent = 1, is also used to set up and expose memory for Request_To_Receive Operations.

Semantics – Request_To_Send (

Op,
Flags,
R_Port,
S_Port,
Key,
S_id,
T_len)

Op = x'06'

Flags (see 6.2) shall specify the Concatenate, Persistent, and Data Channel assignment flags. Concatenate or Persistent shall only = 1 if the corresponding flag was set = 1 by the other end during the Virtual Connection **set up** (see 7.1 and 7.2) and the function is desired for this Operation. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

S_id shall be the Originating Source's Transfer identifier (see 4.4.2) used to identify this Transfer. The Final Destination shall use this value as the R_id parameter when replying to the Originating Source concerning this Transfer.

T_len shall specify the total number of data payload bytes in the Transfer, **or that the size of the Transfer is unlimited** (see 4.4.3).

Issued – By the Originating Source after a Virtual Connection has been established.

Effect – If accepted, the Final Destination shall **set up** to receive the data and then reply back to the Originating Source with the associated parameters. If rejected, the Final Destination shall reply with Reject = 1 in a State_Response Operation (see 4.5.3).

8.2 Request_To_Send_Response

Request_To_Send_Response shall inform the Originating Source whether the **Transfer** was accepted or not. If accepted, the Request_To_Send_Response specifies the Transfer identifier (see 4.4.2) assigned by this end (i.e., the Final Destination) for this Transfer and the number of STUs per Block (see 4.4.6). A Request_To_Send_Response does not give the Originating Source permission to start sending; that comes from a Clear_To_Send. A Clear_To_Send may be used instead of a Request_To_Send_Response if the Final

Destination is able to immediately accept the data.

Semantics – Request_To_Send_Response (

Op,
Flags,
Blocksize [S_count],
R_Port,
S_Port,
Key,
R_id,
S_id,
T_len)

Op = x'07'

Flags (see 6.2) shall specify the Concatenate flag. Concatenate shall only = 1 if Concatenate was set = 1 by the remote end device during the Virtual Connection **set up** (see 7.1 and 7.2). The appropriate value for the Reject and Interrupt flags shall also be carried (see 4.5.3 and 4.5.5).

Blocksize, carried in the S_count field, shall specify the number of STUs in a Block (see 4.4.6).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the Originating Source in the Request_To_Send Operation.

S_id shall be the Transfer identifier (see 4.4.2) used by the Final Destination to identify this Transfer. The Originating Source shall use this value as the R_id parameter when replying to the Final Destination concerning this Transfer.

T_len shall be the same value as in the Request_To_Send Operation, but it need not be checked by the Originating Source (see 4.4.3).

Issued – By the Final Destination.

Effect – The Originating Source shall segment the Transfer into Blocks and STUs for transmission.

8.3 Request_To_Receive

Request_To_Receive, **issued by the Final Destination**, asks for a single Block of data to be sent from a previously allocated location in the Originating Source. The Request_To_Receive specifies the number of data bytes to be sent

from the Originating Source to the Final Destination and whether 64-bit addresses or Buffer Indexes are used at the Originating Source, at the Final Destination, or at both. A Request_To_Receive transfers a single Block; there is no notion of a multi-Block Request_To_Receive data movement.

When a Request_To_Receive is issued, it is assumed that the ULPs on both end devices had previously allocated resources for the entire Transfer through a previous Request_To_Send Operation. Note that the device at either end of the Virtual Connection may issue a Request_To_Receive.

Semantics – Request_To_Receive (

Op,
Flags,
R_Port,
S_Port,
Key,
R_id,
S_id,
Bufx,
Offset,
T_len,
OS_Bufx,
OS_Offset)

Op = x'08'

Flags (see 6.2) shall specify the Concatenate and Source_Concatenate flags. Source_Concatenate shall only = 1 if the corresponding flag was set by the other end device during the Virtual Connection **set up** (see 7.1 and 7.2), and the function is desired for this Operation. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the Originating Source in the Request_To_Send Operation.

S_id shall be the Transfer identifier (see 4.4.2) used by the Final Destination to identify this Transfer. The Originating Source shall use this value as the R_id parameter when replying to the Final Destination concerning this Transfer.

Bufx shall specify the initial Buffer Index in the Final Destination where the data will be placed (see 4.4.7).

Offset is a value that the Final Destination **must** receive with the first STU of the Block so that

the data can be properly placed in the Final Destination's memory (see 4.4.7).

T_len shall specify the total number of data payload bytes in the Transfer (see 4.4.3).

OS_Bufx shall specify the Originating Source's Buffer Index (see 4.4.8).

Open Issue – How does the device issuing the Request_To_Receive know the value to stick in the OS_Bufx parameter?

OS_Offset shall specify the Originating Source's offset value (see 4.4.8).

Issued – By the Final Destination.

Effect – If accepted, the Originating Source shall send the specified Block of data. If rejected, the Originating Source shall reply with Reject = 1 in a State_Response Operation (see 4.5.3).

8.4 Clear_To_Send

Clear_To_Send shall be used to give the Originating Source permission to send one Block. Clear_To_Send may also be used to request retransmission of a Block from systems that are capable of retransmission.

Semantics – Clear_To_Send (

- Op,
- Flags,
- Blocksize [S_count],
- R_Port,
- S_Port,
- Key,
- R_id,
- S_id,
- Bufx,
- Offset,
- T_len,
- B_num,
- I_Offset [OS_Offset])

Op = x'09'

Flags (see 6.2) shall specify the Concatenate and Source_Concatenate flags. These flags shall be the same value as in the Request_To_Send Operation that initiated this Transfer. The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

Blocksize, carried in the S_count field, shall

specify the number of STUs in a Block (see 4.4.6).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the remote end (the Originating Source) of the Virtual Connection.

S_id shall be the Transfer identifier (see 4.4.2) assigned by this end (the Final Destination) of the Virtual Connection.

Bufx shall specify the initial Buffer Index in the Final Destination where the data will be placed (see 4.4.7).

Offset is a value that the Final Destination **must** receive with the first STU of a Block so that the data can be properly placed in the Final Destination's memory (see 4.4.7).

T_len shall be the same value as in the Request_To_Send Operation that this Clear_To_Send is responding to, but it need not be checked by the Originating Source.

B_num shall be the Block number being given permission to be transmitted (see 4.4.4).

I_Offset, carried in the OS_Offset field, shall specify the Offset value associated with the first Block of the Transfer (see 4.4.7).

Issued – By the Final Destination.

Effect – The Originating Source shall send the specified Block.

8.5 Data

A **Data** Operation sends an STU of a Block from the Originating Source to the Final Destination. No STU shall be larger than the maximum STU size determined during the Virtual Connection **set up** (see 7.2).

Semantics – **Data** (
 Op,
 Flags,
 S_count,
 R_Port,
 S_Port,
 Key,
 R_id,
 S_id,
 Bufx,
 Offset,
 Sync [T_len],
 B_num,
 Opaque [OS_Bufx],
 Opaque [OS_Offset])

Op = x'0A'

Flags (see 6.2) shall specify the Interrupt, **Silent**, Send_State, Concatenate, **Last**, and Data Channel assignment flags (see 6.2). Concatenate shall be the same value as in the Request_To_Send or Request_To_Receive that initiated this Transfer. Send_State may be sent with any STU of a Block. The State_Response Control Operation associated with this request shall be issued after processing this STU when Send_State = 1. The sender shall copy (in this **Data** Operation) the Data Channel assignment flags supplied in the corresponding Request_To_Send Operation; the value is a do not care at the receiver.

S_count shall be the STU **number** (see 4.4.6).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the other end (the Final Destination) of the Virtual Connection.

S_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. Note that if this is the first STU associated with a Request_To_Receive Operation, then this Transfer identifier (see 4.4.2) is being assigned by the Originating Source and shall be used by the Final Destination as the R_id parameter when replying to the Originating Source

concerning this Transfer.

Bufx shall be the Buffer Index at the Final Destination (see 4.4.7).

Offset shall be the Final Destination's offset within a Bufx (see 4.4.7).

Sync, carried in the T_len field, shall be a value assigned by the Originating Source to synchronize the current view of the number of empty Slots in the Final Destination (see 4.3.6).

B_num shall be the number of the Block that this STU is a part of.

Opaque data, carried in the OS_Bufx and OS_Offset fields, may be any value assigned by the Originating Source's ULP. The opaque data shall be passed unmodified to the Final Destination's ULP.

Issued – By the Originating Source.

Effect – The Final Destination shall place the STU data in the memory area pointed to by Bufx and offset by the Offset value. The Final Destination shall only accept data into pre-allocated buffer regions. The Final Destination is responsible for ensuring that all of the Blocks of a Transfer are received. The actions to be taken if a Block is missing are beyond the scope of this standard.

8.6 Request_State

Request_State is used to request that the remote end device provide its current number of empty Slots for Schedule Headers, the Block number associated with the last set of contiguously good data received, and whether the named Block was received correctly.

Semantics – Request_State (
 Op,
 Flags,
 R_Port,
 S_Port,
 Key,
 R_id,
 S_id,
 Sync [T_len],
 B_num)

Op = x'0B'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the remote end device of this Virtual Connection. R_id = x'0000' means that the receiver shall not look for a current Transfer and only return the current number of empty Slots for this Virtual Connection.

S_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. If R_id = x'0000', then S_id shall also be x'0000'.

Sync, carried in the T_len field, shall be a value assigned by the local end device (sender) to synchronize the current view of the number of empty Slots in the remote end device (receiver). (See 4.3.6.)

B_num shall indicate the Block number being queried. B_num = x'FFFFFFFF' indicates that the sender does not care about the status of any particular Block.

Issued – By an end device that needs state information from the remote end device of the Virtual Connection. The sender may not have received the State_Response that it expected from a Data Operation and can send a Request_State to recover from a lost or damaged State_Response.

Effect – The receiver shall reply with a State_Response.

8.7 State_Response

State_Response shall be used to indicate the number of empty Slots in this Port of the Virtual Connection (see 4.3.6). State_Response with Reject = 1 is also the response to a Request_To_Receive that was rejected (see 4.5.3). State_Response may also indicate the highest numbered contiguous Block received correctly and whether the Block indicated in the B_num parameter was received correctly (see 4.5.2).

Semantics – State_Response (

Op,
Flags,
C-Slots [S_count],
R_Port,
S_Port,
Key,
R_id,
S_id,
B_seq [Offset],
Sync [T_len],
B_num)

Op = x'0C'

Flags (see 6.2) shall specify Reject = 1 if the immediately previous Request_To_Receive was rejected (see 4.5.3). The appropriate value for the Interrupt flag shall also be carried (see 4.5.5).

C-Slots, carried in the S_count field, shall indicate the sender's view of the number of empty Slots it has available for additional Operations. (See 8.5.) C-Slots = x'FFFF' (i.e., -1) shall indicate that this end device does not implement the Slots mechanism for Operations flow control.

R_id shall echo the S_id value in the Request_State or Data Operation that triggered this State_Response.

S_id shall echo the R_id value in the Request_State or Data Operation that triggered this State_Response. S_id = x'0000' shall mean that the B_seq and B_num parameters are meaningless.

B_seq, carried in the Offset field, shall indicate the highest numbered contiguous Block received correctly. B_seq = x'FFFFFFFF' shall indicate that no Transfers are in progress or no Blocks have been received.

Sync, carried in the T_len field, is echoed from the Request_State, or Data Operation with Send_State = 1, that initiated this State_Response Operation (see 4.3.6).

B_num shall echo the Block number, carried in the B_num field of the Data Operation or the Request_State Operation, if the indicated Block was received correctly. If the indicated Block has not been correctly received, then B_num shall contain x'FFFFFFFF'. The Sync value can be used by the receiving end to identify the Operation containing the B_num being queried.

Table 2 – Virtual Connection Operations summary between end devices A and B

	Op	Flags	S_count	R_Port	S_Port	Key	Bufx	Offset	T_len	B_num
RQP	x'01'	<i>OIUCP</i>	<i>A-Slots</i>	<i>B-Port</i>	<i>A-Port</i>	*	<i>A-BuFSIZE</i>	<i>A-Key</i>	<i>A-Max-STU</i>	<i>EtherType</i>
RQPR	x'02'	<i>OIUCPR</i>	<i>B-Slots</i>	A-Port	<i>B-Port</i>	A-Key	<i>B-BuFSIZE</i>	<i>B-Key</i>	<i>B-Max-STU</i>	*
PT	x'03'	<i>I</i>	*	R_Port	S_Port	R-Key	*	*		*
PTA	x'04'	<i>I</i>	*	R_Port	S_Port	R-Key	*	*		*
PTC	x'05'	<i>I</i>	*	R_Port	S_Port	R-Key	*	*		*

NOTES –

1 – Operation abbreviations:

PT = Port_Teardown

PTA = Port_Teardown_ACK

PTC = Port_Teardown_Complete

RQP = Request_Port

RQPR = Request_Port_Response

2 – Flag abbreviations are: O = Out_of_Order, I = Interrupt, U = Source_Concatenate, C = Concatenate, P = Persistent, R = Reject

3 – R-Key = Key value the receiver binds to, e.g., R-Key = A-Key when Operation issued by device B.

4 – R_Port = Port number in device receiving the Operation, e.g., R_Port = A-Port when issued by device B.

5 – S_Port = Port number in device sending the Operation, e.g., S_Port = B-Port when issued by device B.

6 – The Schedule Header **fields** that are not shown shall be transmitted as zeros.**SYMBOLS –**

* = Unused value, transmit as 0

Values in bold italics are assigned by the specific Operation, and may be used by later Operations

Issued – It is intended that a State_Response be issued by an end device's ULP after receiving Send_State = 1 in a **Data** Operation, or after receiving a Request_State Operation, or to reject an Operation.

Effect – State information is passed to the other end of the Virtual Connection.

8.8 END

END allows either end of the Virtual Connection to terminate a Scheduled Transfer before it has completed and to terminate a Scheduled Transfer of unlimited size.

Semantics – END (

Op,
Flags,
R_Port,
S_Port,
Key,
R_id
S_id)

Op = x'0D'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the other end of the Virtual Connection.

S_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. This S_id value shall not be reused until an END_ACK is received.

Issued – By the Originating Source or the Final Destination.

Effect – A Final Destination receiving an END shall stop sending Control Operations associated with this Scheduled Transfer. An Originating Source receiving an END shall stop sending Control Operations and STUs associated with this Scheduled Transfer. An END kills a Scheduled Transfer, but shall not affect the Virtual Connection carrying the Scheduled Transfer.

Table 3 – Data transfer and status Operations summary between end devices **S and **R****

	Op	Flags	S_count	R_id	S_id	Bufx	Offset	T_len	B_num	OS_Bufx	OS_Offset
RTS	x'06'	ICPD	*	*	S_id	*	*	T_len	*	*	*
RTSR	x'07'	ICR	Blocksize	R_id	S_id	*	*	T_len	*	*	*
RTR	x'08'	IUC	*	R_id	S_id	Bufx	Offset	T_len	*	OS_Bufx	OS_Offset
CTS	x'09'	IUC	Blocksize	R_id	S_id	Bufx	Offset	T_len	B_num	*	I_Offset
Data	x'0A'	ITSCLD	S_count	R_id	S_id	Bufx	Offset	Sync	B_num	Opaque	Opaque
RS	x'0B'	I	*	R_id	S_id	*	*	Sync	B_num	*	*
SR	x'0C'	IR	C-Slots	R_id	S_id	*	B_seq	Sync	B_num	*	*
END	x'0D'	I	*	R_id	S_id	*	*	*	*	*	*
ENDA	x'0E'	I	*	R_id	S_id	*	*	*	*	*	*

NOTES –

1 – Operation abbreviations:

CTS = Clear_To_Send

ENDA = END_ACK

RS = Request_State

RTR = Request_To_Receive

RTS = Request_To_Send

RTSR = Request_To_Send_Response

SR = State_Response

2 – Flag abbreviations are: I = Interrupt, **T = Silent**, S = Send_State, U = Source_Concatenate,C = Concatenate, P = Persistent, **L = Last** STU of Block, R = Reject, D = Data Channel assignment

3 – R_id = Transfer identifier in device receiving the Operation, e.g., R_id = G_id when issued by device H.

4 – S_id = Transfer identifier in device sending the Operation, e.g., S_id = H_id when issued by device H.

5 – Schedule Header parameters that shall be transmitted as assigned in RQP and RQPR Operations:

R_Port = Port number of the device receiving the Operation

S_Port = Port number of the device sending the Operation

Key = Key value assigned by the device receiving the Operation

SYMBOLS -

* = Unused value, transmit as 0

Values in bold italics are assigned by the specific Operation, and may be used by later Operations

8.9 END_ACK

END_ACK confirms that the sending end device has seen, and acted on, the END.

Semantics – END_ACK (

Op,

Flags,

R_Port,

S_Port,

Key,

R_id,

S_id)

Op = x'0E'

Flags shall contain the appropriate value for the Interrupt flag (see 4.5.5).

R_id shall be the Transfer identifier (see 4.4.2) assigned by the remote end device of this Virtual Connection.

S_id shall be the Transfer identifier (see 4.4.2) assigned by this end of the Virtual Connection. This S_id value should not be immediately reused to avoid aliasing.

Issued – By the end of the Virtual Connection that received the END Operation.

Effect – Acknowledgment that the Scheduled Transfer has been terminated.

9 Error processing

Table 5 is a summary of the logged errors. The logging is on a per-Port basis, and shall be available to the ULP that owns the Port. The nature and size of the logs are system dependent.

9.1 Operation timeout

Errors other than syntactic errors are manifested as missing Operations, occurring when the underlying physical media **discard** or **damage** a transmission (see 4.5.4). Such errors are detected by **Op_timeout**, which is system and/or Port dependent. **Op_timeout Occurances shall be logged.** Example means for determining the **Op_timeout** value for a Virtual Connection include:

- a time longer than the measured round-trip time through the software path (use a Request_State / State_Response pair to measure on a per-Port basis); or
- a long fixed time period.

Another system **and/or Port** dependent parameter, Max_Retry, specifies the maximum number of times to retry an Operation. If Max_Retry is reached without success, then the Operation is considered to be aborted and control shall be passed to ULPs. **Max_Retry Occurances shall be logged.**

9.2 Operation Pairs

Table 4 lists the Operation pairs – command and response, or response and completion – that shall be retried if the associated response is not received within an **Op_timeout**.

In addition to the entries in table 4, State_Response is a corresponding pair for **Data** Operations which have Send_State = 1. If the State_Response is not received, then the Originating Source may send a Request_State to obtain the state information.

The ULP in the Final Destination that issues a **Clear_To_Send**, or **Request_To_Receive**, is responsible for timing out these Operations. **The ULP may** or may not use **Op_timeout** to indicate failure.

Table 4 – Operation pairs guarded by **Op_timeout**

Operation	Response
Request_Port	Request_Port_Response
Port_Teardown	Port_Teardown_ACK
Port_Teardown_ACK	Port_Teardown_Complete
Request_To_Send	Request_To_Send_Response, or Clear_To_Send
Request_To_Receive	Data or State_Response
Request_State	State_Response
END	END_ACK

9.3 Syntax errors

9.3.1 Undefined Opcode

An undefined Opcode value may occur due to bit errors, or if the sending device is using a future superset of the Scheduled Transfer Operations. The Operation shall be discarded, an Undefined_Opcode_Error shall be logged, and the Opcode logged in Undefined_Opcode_Value.

9.3.2 Unexpected Opcode

Most of the Operations require previous Operations to **set up** state on each device. If a device receives an out of sequence Opcode (e.g., receiving a Request_Port_Response without sending the initiating **Request_Port**), the Operation shall be discarded, an Unexpected_Opcode_Error shall be logged, and the Opcode logged in Unexpected_Opcode_Value.

9.4 Virtual Connection errors

9.4.1 Invalid Key or Port

All Operations, excluding Request_Port, should have a Key value that validates the Operation for the Virtual Connection. Operations with an invalid Key shall not be executed, and an Invalid_Key_Error shall be logged.

All Operations should have a valid Destination Port value. Operations with an invalid Destination Port value shall not be executed, and an Invalid_Port_Error shall be logged.

NOTE – Multiple contiguous invalid Key and/or Port values may indicate a problem with the link or a malicious host on the network. The supervising process should be informed.

9.4.2 Slots exceeded

Operations that exceed the number of Slots for the Virtual Connection may not be executed, and a Slots_Exceeded_Error shall be logged.

9.4.3 Unknown EtherType

If a Request_Port Operation contains an unknown EtherType, the receiver shall issue a Request_Port_Response with the Reject bit set and log an Unknown_EtherType_Error.

9.4.4 Illegal BuFSIZE

If a Request_Port contains a BuFSIZE value that is less than 256 bytes, then the receiver shall respond with a Request_Port_Response with Reject = 1. If a Request_Port_Response contains a BuFSIZE value that is less than 256 bytes, then the receiver shall respond with a Port_Teardown. In either case, an Illegal_BuFSIZE_Error shall be logged.

9.4.5 Illegal STU size

The maximum STU sizes (A-Max-STU and B-Max-STU) were determined during the Virtual Connection set up (see 4.3.5, 7.1 and 7.2). If the received STU is greater than the maximum STU size, then the STU shall be discarded and an Illegal_STU_Size_Error shall be logged.

9.5 Scheduled Transfer errors

9.5.1 Invalid S_id

All Scheduled Transfer Operations, except Request_To_Send should have a valid Destination id (R_id) for quickly accessing state information for this Scheduled Transfer. After checking the R_id, the S_id should match the stored value for this Transfer. An invalid S_id

shall result in not executing the Operation and logging an Invalid_S_id_Error.

9.5.2 Bad Data Channel specification

During a Request_To_Send Operation, the sending device declares the lower layer Data Channel that will carry Data Operations for this Scheduled Transfer. Some Data Channels may not be available for Scheduled Transfers depending on the lower layer (e.g., b'00' is not a valid choice on HIPPI-6400 as it indicates VC0 which is reserved for Control Operations). The receiver shall issue a Request_To_Send_Response with the Reject bit set.

9.5.3 Concatenate not available

If the Virtual Connection did not specify the capability for Concatenate (see 4.3.7) during the Virtual Connection establishment (see 7.1 and 7.2), any Scheduled Transfer Operations on this Virtual Connection with the Concatenate bit set shall not be executed. The Operation shall be rejected (see 4.5.3) and a Concatenate_Error logged.

9.5.4 Source_Concatenate not available

If the Virtual Connection did not specify the capability for Source_Concatenate (see 4.3.8) during the Virtual Connection establishment (see 7.1 and 7.2), any Scheduled Transfer Operations on this Virtual Connection with the Source_Concatenate bit set shall not be executed. The Operation shall be discarded and a Source_Concatenate_Error logged.

Open Issue – Should we make a more positive response by signaling a Reject?

9.5.5 Persistent not available

If the Virtual Connection did not specify the capability for Persistent (see 4.3.9) during the Virtual Connection establishment (see 7.1 and 7.2), any Scheduled Transfer Operations on this Virtual Connection with the Persistent bit set shall not be executed. The Operation shall be discarded and a Persistent_Error logged.

Open Issue – Should we make a more positive response by signaling a Reject?

9.5.6 Out of Range B_num, Bufx, Offset, S_count, or Sync

During the Clear_To_Send, Data, and State_Response Operations, a Block number may appear that is outside the calculated number of Blocks for the Transfer. If an out of range Block number is encountered, the receiver shall not execute the Operation and shall log an Out_Of_Range_B_num_Error.

If a Data Operation contains a Bufx and/or Offset that exceeds the buffer range allocated by the Final Destination for outstanding Clear_To_Sends, then the receiver shall not execute the Operation and shall log an Out_Of_Range_Bufx_Error.

If a Data Operation contains an Offset larger than the buffer size, the receiver shall not execute the Operation and shall log an Oversized_Offset_Error.

If a Data Operation contains an S_count that is not one greater than the previous STU (for this Block), then the STU is out of order. The receiver shall discard the STU and log an Out_Of_Order_STU_Error.

9.5.7 Request_To_Receive problem

The Request_To_Receive Operation must be set up by a previous Request_To_Send Operation (see 8.3). If an associated Request_To_Send has not been received, then the Request_To_Receive Operation shall be discarded, a State_Response with Reject = 1 sent to the remote end device in response, and a Request_To_Receive_Error logged.

9.5.8 Undefined Flag

If a received Operation contains a flag =1 and use of that flag is not defined for that Operation, then the flag shall be ignored and an Improper_Flag_Use_Error logged.

Open Issue – How long should Sources save their transmit buffers for possible retransmission?

Table 5 – Summary of logged errors

Name	Occurs in Operation
Concatenate_Error	RTS, RTSR, RTR, CTS
Illegal_STU_Size_Error	Data
Illegal_Bufsize_Error	RQP, RQPR
Improper_Flag_Use_Error	all
Invalid_Key_Error	all except RQP
Invalid_Port_Error	all
Invalid_S_id_Error	all with an R_id
Max_Retry_Occurance	END, PT, PTA, RQP, RS, RTR, RTS
Op_timeout_Occurance	END, PT, PTA, RQP, RS, RTR, RTS
Out_Of_Order_STU_Error	Data
Out_Of_Range_B_num_Error	CTS, Data, RS, SR
Out_Of_Range_Bufx_Error	Data
Oversized_Offset_Error	Data
Persistent_Error	RTS
Request_To_Receive_Error	RTR
Slots_Exceeded_Error	all with Opcode ≥ 6
Source_Concatenate_Error	RTR, CTS
Undefined_Opcode_Error	not applicable
Undefined_Opcode_Value	not applicable
Unexpected_Opcode_Error	all except RQP
Unexpected_Opcode_Value	all except RQP
Unknown_EtherType_Error	RQP
Operation abbreviations: CTS = Clear_To_Send PT = Port_Teardown PTA = Port_Teardown_ACK RQP = Request_Port RQPR = Request_Port_Response RS = Request_State RTR = Request_To_Receive RTS = Request_To_Send SR = State_Response	

Annex A (normative)

Using lower layer protocols

A.1 HIPPI-6400-PH as the lower layer

ANSI X3.xxx defines HIPPI-6400-PH, portions of which are repeated here as an aid to the reader. As shown in figure A.1, HIPPI-ST Operations shall be carried over HIPPI-6400-PH with the first eight bytes of the HIPPI-ST Schedule Header occupying the last eight bytes of the HIPPI-6400-PH Header micropacket.

All HIPPI-ST Control Operations shall be carried on HIPPI-6400-PH Virtual Channel VC0. Data Operations shall use Virtual Channel 1, 2, or 3 as specified in the HIPPI-ST Data Channel Assignment flag bits (see 6.2) and carried in a Request_To_Send Operation (see 8.1).

HIPPI-ST shall also specify the EtherType value that is placed in the HIPPI-6400-PH MAC Header (see the reference for RFC 1700 in 4.3.2).

M_len (in the HIPPI-6400-PH MAC Header), specifies the number of bytes following M_len, exclusive of any padding in the last micropacket. Hence, M_len will have the following values:

- M_len = 48 for Control Operations without an optional payload (i.e., 48 = 8 byte IEEE 802.2 LLC/SNAP Header + 40-byte HIPPI-ST Schedule Header);
- M_len = 80 for Control Operations with optional payload;
- M_len = (48 + number of user data payload bytes) for Data Operations.

A.2 HIPPI-FP as the lower layer

ANSI X3.210 defines HIPPI-FP, portions of which are repeated here as an aid to the reader. As shown in figure A.2, HIPPI-ST Operations shall be carried over HIPPI-FP in the D2_Area. The HIPPI-FP D1_Area shall not be used. The HIPPI-FP D2_Offset shall be set to zero. Short bursts shall only be used at the end of a packet, i.e., short first burst is disallowed. Note that $D2_Size = M_len + 16$.

The HIPPI-6400-PH MAC and LLC/SNAP Headers are defined in ANSI X3.xxx, portions of which are repeated here as an aid to the reader. The MAC and LLC/SNAP headers are included to facilitate translation to other protocols. The 48-bit ULA addresses allow address assignment and usage common to other networking technologies. The 12 least significant bits of the 48-bit S_ULA (Source) and D_ULA (Destination) addresses shall be identical to the corresponding addresses carried in the HIPPI Switch Control I-Field (see ANSI X3.222, HIPPI-SC).

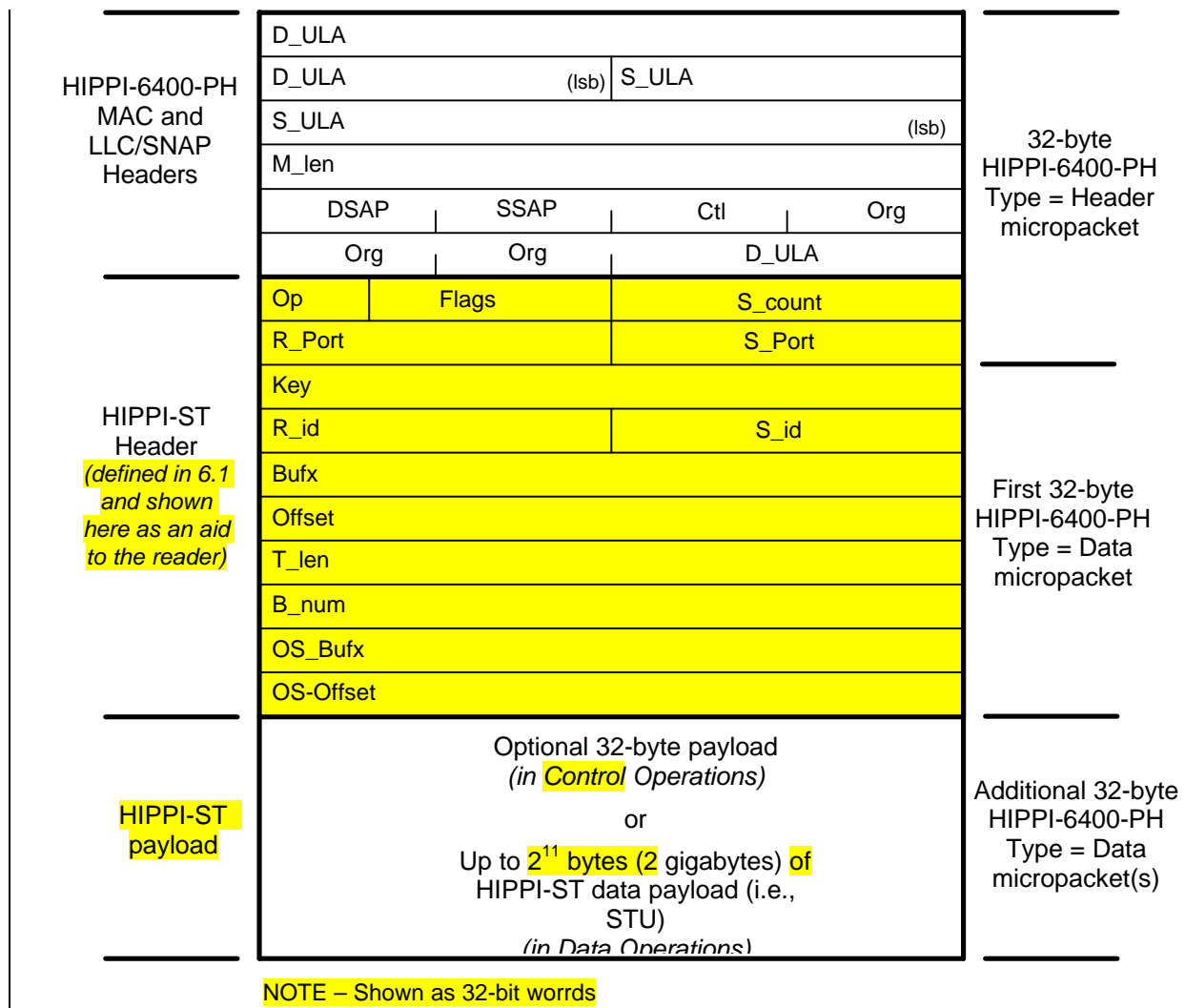


Figure A.1 – HIPPI-ST Operations carried in HIPPI-6400-PH Messages

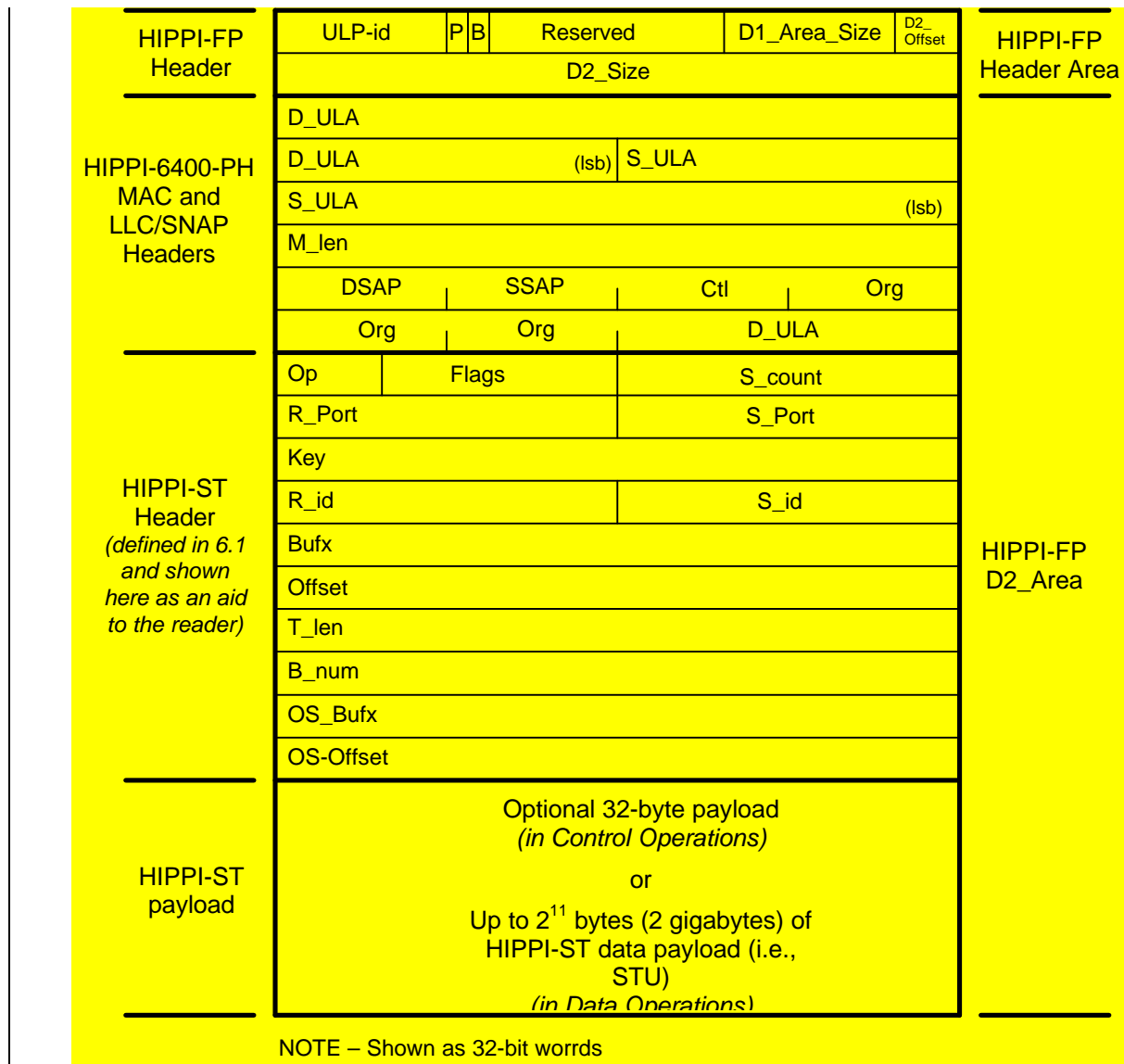


Figure A.2 – HIPPI-ST Operations carried in HIPPI-FP packets

Annex B (informative)

HIPPI-ST striping

Open Issue - New text and figures, complements of Roger Ronald. It needs review.

B.1 Striping principles

HIPPI-ST is capable of supporting multiple physical interfaces for a single Transfer (see figures B.1–B.3). This striping capability may be of benefit when a single interface is not able to support required data rates. It may be especially useful where data is moved from many slower interfaces to a single faster interface or vice-versa. It may also be used with multiple interface at both the Originating Source and Final Destination.

Each STU contains sufficient information to completely identify a particular individual Transfer as well as which piece of the Transfer is carried within. Thus, there is no problem in allowing the Originating Source location to change during a Transfer, so long as the STU arrives intact at the Final Destination. In the other direction of movement, the Final Destination may be specified as multiple physical locations via a differing MAC Source ULA in CTS operations.

HIPPI-ST supports striping without any striping specific protocol constructs. The only difference between striped and non-striped operation is the selection of port MAC addresses to allow concurrent data movement. Striping should always be performed on a Block, not a STU basis. Otherwise, STU in-order delivery is unlikely to be guaranteed.

There are a few conventions that should be followed to facilitate striping:

- Block sizes (when striping is desirable) must be small enough to support concurrency and allow each channel to have at least one Block to send.
- Sufficient CTS operations should be kept outstanding by a data receiver to allow concurrent Data operations.
- The interface adapter(s) must be capable of handling multiple Blocks simultaneously. This

may require communication between interfaces (or their software drivers) within a system.

- The return ULA for each operation is specified by the Source ULA for that operation. HIPPI-ST implementers should not assume that the Source ULA for a given port will remain constant.

B.2 Many-to-one striping

When a number of lower bandwidth interfaces are aggregated to communicate with one faster interface (using a translator or bridge), striping the lower bandwidth interfaces together can allow legacy systems to communicate quickly over newer network infrastructures. In this case, action to implement striping is required only on the side of the lower bandwidth interface.

After port set-up, a Transfer is initiated with a RTS operation. An RTS_Response will be received, either as a discrete message or as part of a CTS. As CTS operations are received, the system with multiple lower bandwidth ports can move a Block of data for each CTS received. As many Blocks can be in transit concurrently as there are ports to send them on and CTS operations authorizing them.

The system receiving these Blocks processes them normally, placing them into memory as their Bufx and Offset values dictate.

B.3 One-to-many striping

Transfers made from one high-speed interface can also be spread across more than one low-speed interface without any special action on the part of the high speed system.

After port set up, the Transfer is initiated with a RTS operation from the high-speed interface. The low-speed interface that has done the port set up will return a RTS_Response, either as a

discrete message or as part of a CTS. Each CTS sent should be sent from the channel where it is desired that the data be received. An alternative is for all CTS operations to be sent from the same interface, "spoofing" the Source ULA for the operation to make it appear that it was generated by the interface where the data is desired. Using this "spoofing" substitution in combination with a dedicated control channel may also prevent or reduce blocking effects where the underlying physical media suffers from high latency. Subsequent Data Operations may then be done concurrently and will use the Source ULA from in the respective

Clear_To_Send Operation as the Destination ULA.

B.4 Many-to-many striping

Many-to-many striping is the combination of the one-to-many and many-to-one striping. The system receiving data indicates his desire to receive in a striped fashion by issuing multiple CTS operations with differing return Source ULAs. The system sending data chooses to stripe by sending from multiple interfaces that are capable of reaching the proper destination ULA.

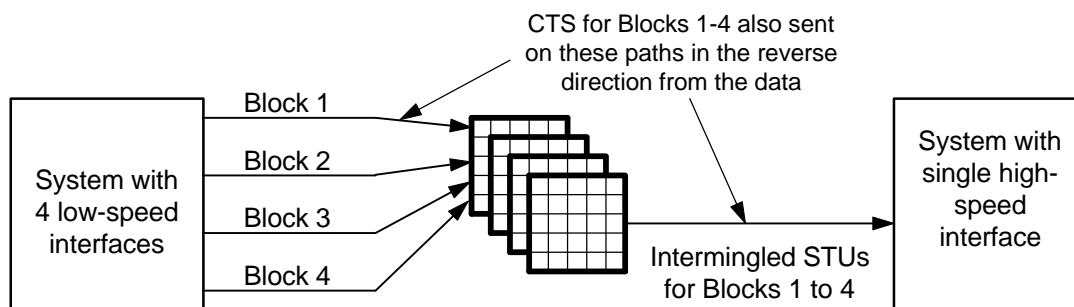


Figure B.1 – Many-to-one striping

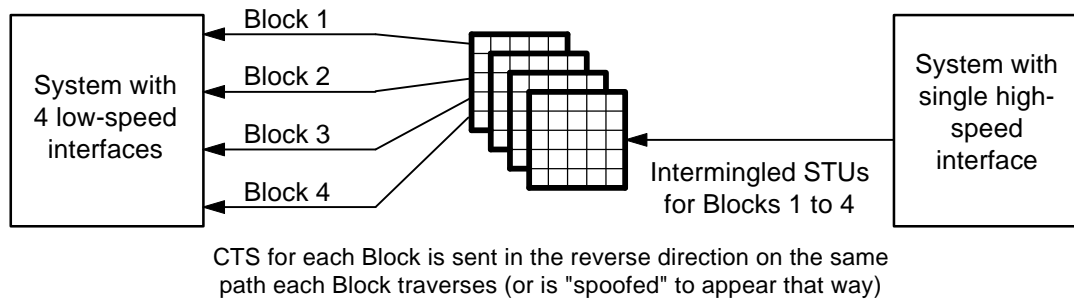


Figure B.2 – One-to-many striping

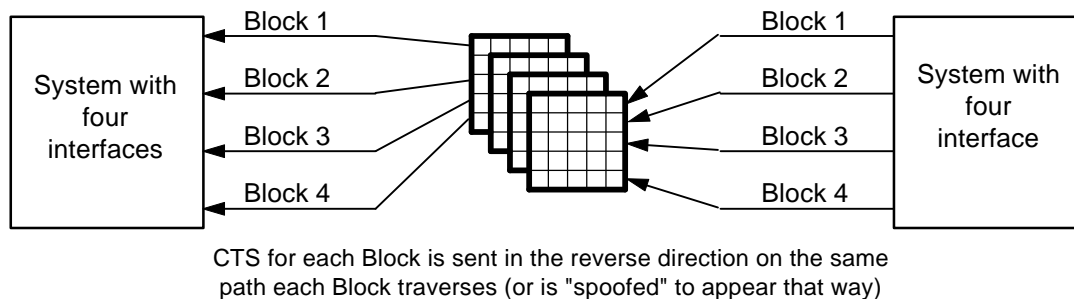


Figure B.3 – Many-to-many striping

Annex C
(informative)

Scheduled Transfer example

Open Issue – This annex has not been updated to match the rest of the document. Hence, the original text has been deleted as being more confusing than helpful. It will be updated and added again as time permits.